



การพัฒนาเวิร์กโฟลว์และเว็บเซอร์วิสที่เหมาะสมในการวิเคราะห์ SNP
Development of Suitable Workflows and Web Services for SNP Analysis

กษิดิ์กฤษณ์ คำเกลี้ยง

Kasikrit Damkliang

วิทยานิพนธ์นี้สำหรับการศึกษาตามหลักสูตรปริญญา
วิศวกรรมศาสตรมหาบัณฑิต สาขาวิชาวิศวกรรมคอมพิวเตอร์
มหาวิทยาลัยสงขลานครินทร์

A Thesis Submitted in Fulfillment of the Requirements for the Degree of
Master of Engineering in Computer Engineering
Prince of Songkla University

2552

ลิขสิทธิ์ของมหาวิทยาลัยสงขลานครินทร์

ชื่อวิทยานิพนธ์ การพัฒนาเวิร์คโฟลว์และเว็บเซอร์วิสที่เหมาะสมในการวิเคราะห์สนิป
ผู้เขียน นายภิชิต์กฤษณ์ คำเกลี้ยง
สาขาวิชา วิศวกรรมคอมพิวเตอร์

อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก

คณะกรรมการสอบ

.....
(ผู้ช่วยศาสตราจารย์ ดร.พิชญา ตันชัยย์)

.....ประธานกรรมการ
(ผู้ช่วยศาสตราจารย์ ดร.วิภาดา เวทย์ประสิทธิ์)

อาจารย์ที่ปรึกษาวิทยานิพนธ์ร่วม

.....กรรมการ
(รองศาสตราจารย์ ดร.วสันต์ จันทราทิตย์)

.....
(ผู้ช่วยศาสตราจารย์ ดร.สุนทร วิฑูรพจน์)

.....กรรมการ
(ผู้ช่วยศาสตราจารย์ ดร.พิชญา ตันชัยย์)

.....กรรมการ
(ผู้ช่วยศาสตราจารย์ ดร.สุนทร วิฑูรพจน์)

บัณฑิตวิทยาลัย มหาวิทยาลัยสงขลานครินทร์ อนุมัติให้รับวิทยานิพนธ์ฉบับนี้
สำหรับการศึกษา ตามหลักสูตรปริญญาวิศวกรรมศาสตรมหาบัณฑิต สาขาวิชาวิศวกรรม
คอมพิวเตอร์

.....
(รองศาสตราจารย์ ดร.เกริกชัย ทองหนู)
คณบดีบัณฑิตวิทยาลัย

ชื่อวิทยานิพนธ์	การพัฒนาเวิร์คโฟลว์และเว็บเซอร์วิสที่เหมาะสมในการวิเคราะห์สนิป
ผู้เขียน	นายภิชิต์กฤษณ์ คำเกลี้ยง
สาขาวิชา	วิศวกรรมคอมพิวเตอร์
ปีการศึกษา	2552

บทคัดย่อ

วิทยานิพนธ์นี้ นำเสนอการพัฒนาเวิร์คโฟลว์เพื่อวิเคราะห์สนิปของกึ่งและมนุษย์ โดยใช้เทคโนโลยีมายกริด โปรแกรมทาเวอร์นา เว็บเซอร์วิส และนำเสนอกระบวนการที่เหมาะสมในการเลือกเว็บเซอร์วิสและวิธีการพัฒนาหรือการวาดเวิร์คโฟลว์ โดยใช้งานวิจัยการวิเคราะห์สนิป ทั้งสองงานเป็นกรณีศึกษา

เวิร์คโฟลว์ในกรณีศึกษาเรื่องสนิปของกึ่ง มักเกิดปัญหาใช้เวลานานและได้ผลลัพธ์ไม่คงเส้นคงวาเมื่อจำนวนข้อมูลเกินขอบเขตที่บริการภายนอกรับได้ วิทยานิพนธ์นี้จึงนำเสนอการพัฒนาบริการท้องถิ่นขึ้นมาใช้งานเอง เพื่อแก้ไขปัญหาดังกล่าว

ในงานวิเคราะห์สนิปมนุษย์ เวิร์คโฟลว์มีความซับซ้อนและมีหลายชั้น วิทยานิพนธ์นี้นำเสนอกลไกการลดเวลาในการพัฒนาและทดสอบเวิร์คโฟลว์ โดยพัฒนาเวิร์คโฟลว์สำหรับตรวจสอบความถูกต้องของเวิร์คโฟลว์เป้าหมายก่อนการทำงานจริง ว่าทันสมัยพร้อมทำงานหรือล้าสมัยไปแล้วและให้ข้อมูลที่เกี่ยวข้องเพื่อการปรับแก้เวิร์คโฟลว์ ซึ่งสามารถลดเวลาที่เสียไปในการตรวจสอบและหาที่ผิดได้เป็นอย่างมาก

คำสำคัญ: เวิร์คโฟลว์, ทาเวอร์นา, เว็บเซอร์วิส, ไบโอมาร์ท, สนิป, การตรวจสอบความถูกต้อง

Thesis Title	Development of suitable workflows and Web services for SNP analysis
Author	Mr. Kasikrit Damkliang
Major Program	Computer Engineering
Academic Year	2009

ABSTRACT

This thesis presents workflow development for shrimp and human SNP analysis using myGrid technology, Taverna, and Web service technology. Through the two case studies of shrimp and human SNP analysis, this thesis suggests suitable processes in Web service selection and workflow composition.

The first workflow for shrimp SNP analysis sometimes fails due to time-out and inconsistency problems caused by external services when the number of data exceeds the service's limit. Our solution is providing local Web service interfaces instead.

In case of human SNP analysis, the workflow is complex and nested. We propose a time reduction mechanism in development and testing using a validation workflow. This validation workflow is to be used before running a workflow to show whether there are up-to-date or out-of-date. It also provides associate information for re-composing and configuring the workflow. Therefore, it can greatly reduce the debugging and checking times.

Keywords: Workflow, Taverna, Web Service, SNP, Validation

กิตติกรรมประกาศ

ขอขอบพระคุณผู้ช่วยศาสตราจารย์ ดร.พิชญ์ ตันทัยย์ อาจารย์ที่ปรึกษาหลักที่ได้คำปรึกษาและความรู้ในด้านต่างๆ พร้อมทั้งชี้แนะแนวทางการดำเนินการเกี่ยวกับวิทยานิพนธ์ รวมถึงการให้โอกาสและประสบการณ์ในด้านวิชาการต่างๆทั้งในประเทศและต่างประเทศ อีกทั้งสนับสนุนอุปกรณ์ในการทำวิจัย ตลอดจนช่วยตรวจและแก้ไขปรับปรุงวิทยานิพนธ์ให้เป็นอย่างสมบูรณ์ และที่สำคัญคือการมีอาจารย์เป็นกัลยาณมิตรธรรม

ขอขอบพระคุณผู้ช่วยศาสตราจารย์ ดร.สุนทร วิฑูรพจน์ อาจารย์ที่ปรึกษาร่วม ที่ได้คอยชี้แนะแนวทางต่างๆที่เป็นประโยชน์ และรวมทั้งได้ให้ความรู้ความเข้าใจที่จำเป็นต่อการทำวิจัย

ขอขอบพระคุณคณะกรรมการสอบวิทยานิพนธ์ที่ได้ให้คำแนะนำในการปรับปรุงวิทยานิพนธ์ให้สมบูรณ์ยิ่งขึ้น

ขอขอบพระคุณคณาจารย์และบุคลากรในภาควิชาวิศวกรรมคอมพิวเตอร์ทุกท่าน ที่ให้คำปรึกษาและความช่วยเหลือในระหว่างทำวิทยานิพนธ์

ขอขอบพระคุณ คณะวิศวกรรมศาสตร์ บัณฑิตวิทยาลัย มหาวิทยาลัยสงขลานครินทร์ โครงการสร้างโครงสร้างพื้นฐานด้วยไบโอกริดภายใต้การสนับสนุนของศูนย์ไทยกริดแห่งชาติและสำนักงานส่งเสริมอุตสาหกรรมซอฟต์แวร์แห่งชาติ (องค์การมหาชน) ศูนย์กริดแห่งมหาวิทยาลัยสงขลานครินทร์ และศูนย์ความเป็นเลิศด้านชีววิทยาศาสตร์ของประเทศไทย (TCELS) คณะแพทยศาสตร์ โรงพยาบาลรามาธิบดี มหาวิทยาลัยมหิดล ที่ได้มอบทุนสนับสนุนในการทำวิจัย

ขอขอบคุณนักศึกษาปริญญาโทสาขาวิศวกรรมคอมพิวเตอร์ทุกท่านที่ให้คำแนะนำและเป็นกำลังใจมาโดยตลอด

ขอขอบพระคุณ พ่อแม่และทุกคนในครอบครัว ผู้ให้กำลังใจและความหวังในการฝ่าฟันอุปสรรค และสนับสนุนให้การเรียนรู้ได้สำเร็จลุล่วงไปด้วยดี

ท้ายที่สุด ขอกราบขอบพระคุณพระอาจารย์ปราโมทย์ ปาโมชโช และครูบาอาจารย์ทุกท่านทั้งที่ยังดำรงชั้นอยู่และได้ละชั้นไปแล้ว ในคำเทศนาธรรมและคำสอนอันเป็นทางเอก ที่ทำให้ผู้ทำวิทยานิพนธ์ได้ยืนหยัดและยอมรับไปตามความเป็นจริง

กษิติภัณฑ์ คำเกลี้ยง

สารบัญ

	หน้า
สารบัญ.....	(6)
รายการตาราง.....	(11)
รายการภาพประกอบ.....	(13)
ลักษณะคำย่อและตัวย่อ.....	(18)
บทที่	
1 บทนำ.....	1
1.1 ที่มาและความสำคัญของวิทยานิพนธ์.....	1
1.2 การตรวจเอกสาร (Literatures Review).....	5
1.2.1 โปรแกรมทาเวอร์น่า.....	5
1.2.2 โปรแกรมทาเวอร์น่าสอง รุ่นทดลองใช้.....	6
1.2.3 บริการไบโอมาร์ท (BioMart service).....	7
1.2.4 การพัฒนาเวิร์คโฟลว์สำหรับวิเคราะห์สลับของมนุษย์.....	9
1.3 วัตถุประสงค์.....	11
1.4 ขอบเขตของงานวิจัย.....	11
1.5 ประโยชน์ที่คาดว่าจะได้รับ.....	12
1.6 สรุป.....	12
2 สถาปัตยกรรม โปรแกรมทาเวอร์น่าและเวิร์คโฟลว์.....	14
2.1 สถาปัตยกรรมที่เกี่ยวข้องกับวิทยานิพนธ์.....	14
2.1.1 การพัฒนาระบบบนเครือข่ายกริดด้วยแนวความคิดของสถาปัตยกรรมเชิงบริการ.....	14
2.1.2 เว็บเซอร์วิส.....	15
2.1.3 สถาปัตยกรรมของโปรแกรมทาเวอร์น่า.....	17
2.2 โปรแกรมทาเวอร์น่าและเวิร์คโฟลว์.....	19
2.2.1 ภาพรวมของโปรแกรมทาเวอร์น่า.....	19

สารบัญ (ต่อ)

บทที่	หน้า
2.2.2 ส่วนประกอบของโปรแกรมทาเวอร์น่า	22
2.3 Scufi และลักษณะเด่นของโปรแกรมทาเวอร์น่า	33
2.3.1 การทำงานซ้ำๆ โดยปริยาย (Implicit iteration)	33
2.3.2 บีนเชลล์สคริปต์ (Beanshell scripting)	34
2.3.3 บริการ Soaplab	36
2.3.4 การดีบั๊กและชี้แนะทิศทางการทำงานของเวิร์คโฟลว์ด้วยเบรคพอยต์	37
2.4 ทาเวอร์น่าสอง รุ่นทดลองใช้งาน (Taverna 2 Preview)	39
2.5 บริการ ไบโอมาร์ท (BioMart)	42
2.5.1 แนะนำบริการ ไบโอมาร์ท	42
2.5.2 คำอธิบาย	42
2.5.3 สถาปัตยกรรมของไบโอมาร์ท	43
2.5.4 การเข้าถึงข้อมูลของบริการ ไบโอมาร์ท	45
2.6 การคิวรีข้อมูลด้วยบริการ ไบโอมาร์ทในโปรแกรมทาเวอร์น่า	50
2.6.1. การเพิ่มบริการ ไบโอมาร์ท	50
2.6.2. การปรับแต่งแอตทริบิวต์ (Configuring Attributes)	52
2.6.3. โหมดผลลัพธ์ของบริการ ไบโอมาร์ท (Result modes)	53
2.7 สรุป	54
3 การเลือกบริการที่เหมาะสมและการสร้างบริการท้องถิ่นในการวิเคราะห์สนิปคู่	56
3.1 บทนำ	56
3.2 ขั้นตอนการเลือกบริการที่เหมาะสม	57
3.2.1 การค้นหาบริการจากเว็บพอร์ทัลของมายกริด	57
3.2.2 การค้นหาบริการจากอินเทอร์เน็ตโดยปลั๊กอินของโปรแกรมทาเวอร์น่า	58
3.3 ขั้นตอนการพัฒนาเวิร์คโฟลว์	60
3.4 การออกแบบเวิร์คโฟลว์และเว็บเซอร์วิสสำหรับการวิเคราะห์สนิปคู่	64

สารบัญ (ต่อ)

บทที่	หน้า
3.4.1 กรณีศึกษาและเบื้องหลังของงาน	64
3.4.2 วิธีการทำงานแบบเดิม (Copying and Pasting).....	66
3.4.3 บริการหรือเว็บเซอร์วิสที่เลือกใช้งานและการพัฒนาเวิร์คโฟลว์	67
3.4.4 การแก้ไขปัญหาการเดินเวลาและความไม่คงเส้นคงวาในการทำงาน	74
3.4.5 การพัฒนาบริการท้องถิ่น	76
3.4.6 การพัฒนาเวิร์คโฟลว์ที่ใช้บริการท้องถิ่น	91
3.4.7 ผลการทดลอง.....	96
3.5 อภิปรายผลการทดลอง.....	99
3.6 สรุป	100
4 การออกแบบและการพัฒนาเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์	101
4.1 ที่มาและความสำคัญ	101
4.2 ประโยชน์ในด้านการพัฒนาเวิร์คโฟลว์	101
4.3 ขั้นตอนการค้นหาบริการที่เหมาะสม	103
4.4 การออกแบบระบบงานในงานวิจัยด้านสนิปของมนุษย์	103
4.4.1 กระบวนการค้นหาข้อมูล.....	103
4.4.2 Data flow diagram ของเวิร์คโฟลว์.....	105
4.5 การสร้างเวิร์คโฟลว์ของแต่ละกระบวนการ (Implementation)	106
4.6 การรวบรวมทุกกระบวนการเป็นเวิร์คโฟลว์เดียวกัน	117
4.7 การทดลองวัดเวลาการทำงานของเวิร์คโฟลว์	121
4.8 สรุป	123
5 การออกแบบ การพัฒนาและการทดสอบเวิร์คโฟลว์สำหรับตรวจสอบความถูกต้องของบริการ ไปโอมาร์ท	124
5.1 เกริ่นนำ	124

สารบัญ (ต่อ)

บทที่	หน้า
5.2 แนวคิดการออกแบบเวิร์คโฟลว์สำหรับการตรวจสอบบริการไป โอมาร์ท.....	125
5.2.1 โครงสร้างของ Scufi.....	125
5.2.2 โครงสร้างของบริการไป โอมาร์ทที่ใช้งานวิจัย	128
5.3 การออกแบบ Data flow ของเวิร์คโฟลว์สำหรับการตรวจสอบบริการไป โอมาร์ท	129
5.3.1 การทำงานของเวิร์คโฟลว์ในภาพรวม	130
5.3.2 กลไกการตรวจสอบบริการไป โอมาร์ท	132
5.4 การพัฒนาเวิร์คโฟลว์	133
5.4.1 การดึงอิลิเมนต์ข้อมูลในเอกสาร Scufi ของเวิร์คโฟลว์.....	134
5.4.2 การสืบค้นข้อมูลรายชื่อฟิลเตอร์และแอ็คตริวิตีที่ทันสมัย	135
5.4.3 การเปรียบเทียบข้อมูลในเวิร์คโฟลว์กับริจิสทรีของไป โอมาร์ท	136
5.4.4 การปรับแต่งกลไกการทำซ้ำของพอร์ตอินพุตในเวิร์คโฟลว์.....	138
5.5 การตรวจสอบเวิร์คโฟลว์สำหรับการตรวจสอบบริการไป โอมาร์ท	142
5.6 สรุป	146
6 การตรวจสอบบริการไป โอมาร์ทในเวิร์คโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์.....	147
6.1 การตรวจสอบเวิร์คโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์.....	147
6.2 กระบวนการหลังการตรวจสอบ	153
6.2.1 การปรับปรุงและการทดสอบซ้ำเวิร์คโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์.....	153
6.2.2 การตรวจสอบความถูกต้องของผลลัพธ์.....	157
6.3 ข้อเสนอแนะ	157
6.4 สรุป	158
7 สรุป.....	159
7.1 ภาพรวมของงานวิทยานิพนธ์.....	159
7.2 การเลือกบริการที่เหมาะสมกับลักษณะของงานและการพัฒนาเวิร์คโฟลว์.....	160

สารบัญ (ต่อ)

บทที่	หน้า
7.3 บริการท้องถิ่นในการทำนายโครงสร้างต้นสายวิวัฒนาการของกุ้ง	160
7.4 เวิร์คโฟลว์สำหรับวิเคราะห์หัตถ์สปีของมนุษย์	161
7.5 เวิร์คโฟลว์สำหรับการตรวจสอบบริการไปโอมาร์ท.....	161
7.6 อุปสรรคและปัญหา	162
7.7 ข้อเสนอแนะ	163
เอกสารอ้างอิง	165
รายการภาคผนวก	
ภาคผนวก ก	170
ภาคผนวก ข	176
ภาคผนวก ค	182
ประวัติผู้เขียน	183

รายการตาราง

ตาราง.....	หน้า
2.1 ชนิดข้อมูลที่โปรแกรมทาเวอร์นารองรับ.....	25
2.2 เมตาดาตาและรายละเอียดการปรับแต่งของบริการไบโอมาร์ท	48
2.3 พารามิเตอร์ที่ใช้ในการปรับแต่งเมตาดาตา.....	48
3.1 แสดงตัวอย่างการทำงานของบริการ ‘Split string into string list’	69
3.2 บริการ EBI ที่ใช้ในการพัฒนาเวิร์ค โฟลว์วิเคราะห์ของกึ่ง	70
3.3 พารามิเตอร์ของแต่ละบริการหรือเว็บเซอร์วิสที่กำหนดในเวิร์ค โฟลว์.....	71
3.4 กรณีทดสอบและผลการทดสอบ	74
3.5 กรณีทดสอบและผลการทดสอบของเวิร์ค โฟลว์ใหม่ที่ใช้บริการท้องถิ่น	96
3.6 ตารางการเปรียบเทียบกรณีทดสอบและผลการทดสอบของวิธีการทำงานแบบ	97
4.1 ตัวอย่างการจัดการอินพุตและเอาต์พุตในเวิร์ค โฟลว์.....	108
4.2 การทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสท้องถิ่นในการค้นหาข้อมูลยีน	109
4.3 การทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสท้องถิ่นในการค้นหาข้อมูล	
Gene Ontology	110
4.4 ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล OMIM.....	112
4.5 ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล Motif.....	114
4.6 ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล	
Pathways	116
4.7 คุณสมบัติของเวิร์ค โฟลว์ค้นหาข้อมูลสลับและระยะเวลาในการทำงาน.....	120
5.1 XPaht expression ที่ใช้ดึงอีลิเมนต์ของข้อมูลในเวิร์ค โฟลว์.....	134
5.2 รูปแบบ URL สำหรับการสืบค้นรีจิสทรีของบริการไบโอมาร์ท.....	135
5.3 ตัวอย่างของข้อมูลอินพุตแบบหลายรายการ	139

รายการตาราง (ต่อ)

ตาราง.....	หน้า
5.4 ผลลัพธ์การทำให้เข้าแบบ Cross product	140
5.5 ผลลัพธ์การทำให้เข้าแบบ Dot product	140
5.6 ผลลัพธ์กลไกการทำให้เข้าของข้อมูลอินพุตในเวิร์ค โฟลว์	141
5.7 ตัวอย่างการทำงานของบริการท้องถิ่น Java.....	144
6.1 คุณสมบัติและผลการทดสอบเวิร์ค โฟลว์เกสซ์พันธุศาสตร์	153

รายการรูป

รูป	หน้า
1.1 สรุปงานในวิทยานิพนธ์.....	4
1.2 โปรแกรมทาวเวอร์นาแสดงเวิร์คโฟลว์และเว็บเซอร์วิส ดยการแบ่งประเภทด้วยสีต่างๆ	6
1.3 โครงสร้างบริการไปโอมาร์ทในงานวิจัย	8
2.1 การทำงานของสถาปัตยกรรมเชิงบริการ	15
2.2 สถาปัตยกรรมการทำงานของเว็บเซอร์วิส	17
2.3 สถาปัตยกรรมของระบบที่เกี่ยวข้องในงานวิทยานิพนธ์.....	18
2.4 หน้าต่าง AME ของโปรแกรมทาวเวอร์นา.....	23
2.5 รายละเอียดของคำอธิบายใน AME.....	24
2.6 การเชื่อมต่ออินพุตเข้ากับโปรเซสเซอร์	26
2.7 รายละเอียด Remote resource usage ของโปรเซสเซอร์ทั้งหมดที่ใช้ในเวิร์คโฟลว์.....	27
2.8 กำหนดค่าการคงทนต่อการทำงาน	28
2.9 ไดอะแกรมของเวิร์คโฟลว์	29
2.10 Available Service Panel แสดง โหนดของบริการ Soaplab จาก EBI	30
2.11 ผลลัพธ์การทำงานของเวิร์คโฟลว์จาก Result Browser	31
2.12 รายละเอียดการทำงานของรายงานโปรเซส	32
2.13 Configure iterators ในโปรแกรมทาวเวอร์นา.....	34
2.14 เวิร์คโฟลว์การใช้งานบีนเชลล์อย่างง่าย	35
2.15 ตัวอย่างสคริปต์ของบีนเชลล์	36
2.16 สถาปัตยกรรมการทำงานของ Soaplab Analysis Tool	37
2.17 ตัวอย่างกลไกการทำงานของเวิร์คโฟลว์เมื่อมีการกำหนดเอ็นพอยต์.....	39
2.18 เวิร์คโฟลว์ในการค้นหา Gene Ontology จาก Ensembl IDs	40

รายการรูป (ต่อ)

รูป	หน้า
2.19 Health check report ของเวิร์ค โฟลว์	41
2.20 สถานะการทำงานของเวิร์ค โฟลว์	41
2.21 ภาพรวมของไบ โอมาร์ท.....	43
2.22 สถาปัตยกรรมแบบ Three Tier Architecture ของไบ โอมาร์ท	44
2.23 ความสัมพันธ์ระหว่างดาตามาร์ทและดาตาเซ็ต	44
2.24 ตัวอย่างเมตาดาตาที่ใช้ในการคิวรีข้อมูลจากไบ โอมาร์ท	47
2.25 ตัวอย่างขอบเขตของ Virtual schema ใดๆของดาตามาร์ท.....	47
2.26 ตัวอย่างของคำสั่งการคิวรีในดาตาเซ็ต.....	49
2.27 โครงสร้างต้นไม้ของบริการไบ โอมาร์ทในโปรแกรมทาวเวอร์น่า	51
2.28 หน้าต่าง ‘Out of the box’ กำหนดค่าฟิลเตอร์ต่างๆ	52
2.29 หน้าต่าง ‘Out of the box’ กำหนดค่าแอ็ตทริบิวต์	53
2.30 บริการไบ โอมาร์ทที่มีหลายเอาท์พุต.....	54
3.1 เว็บไซต์ Biological web services ของมายกริด	58
3.2 เงื่อนไขต่างๆในการค้นหาบริการของปลั๊กอิน Feta ในโปรแกรมทาวเวอร์น่า	59
3.3 ขั้นตอนการหาบริการที่เหมาะสมและการพัฒนาเวิร์ค โฟลว์ที่เสนอในวิทยานิพนธ์นี้	62
3.4 Data flow diagram อธิบายข้อมูลที่ไหลในกระบวนการพัฒนาเวิร์ค โฟลว์.....	63
3.5 ภาพรวมการติดต่อสื่อสารของโปรแกรมทาวเวอร์น่าและจุดที่เกิดปัญหาการกินเวลา.....	65
3.6 ลำดับของกระบวนการวิเคราะห์สลิปของกุ่ม.....	68
3.7 เวิร์ค โฟลว์แรก que พัฒนาขึ้นโดยใช้บริการต่างๆที่กระจายอยู่บนอินเทอร์เน็ต	72
3.8 ตัวอย่างอินพุตลำดับนิวคลีโอไทด์รูปแบบ Fasta ของกุ่มแซบว้าย	73
3.9 ข้อความแสดงความผิดพลาดในการใช้บริการของ EBI ที่โปรแกรมทาวเวอร์น่ารายงาน.....	74

รายการรูป (ต่อ)

รูป	หน้า
3.10 สถาปัตยกรรมโดยรวมของบริการท้องถิ่น	75
3.11 Flowchart ของกระบวนการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ	77
3.12 โครงสร้างต้นไม้สายวิวัฒนาการในรูปแบบ Treefile ของลำดับนิวคลีโอไทด์	78
3.13 โครงสร้างต้นไม้สายวิวัฒนาการในรูปแบบ Text ของลำดับนิวคลีโอไทด์	78
3.14 สถาปัตยกรรมการทำงานของโปรแกรมต่างๆของบริการท้องถิ่น	79
3.15 Sequence diagram ของการทำงานภายในเวิร์กโฟลว์ของโปรแกรมทาวเวอร์นาในการทำนาย โครงสร้างต้นไม้สายวิวัฒนาการ	81
3.16 เอกสาร ACD สำหรับโปรแกรม fseqboot	82
3.17 เอกสาร ACD สำหรับโปรแกรม fdnpars	83
3.18 เอกสาร ACD สำหรับโปรแกรม fprotpars	84
3.19 เอกสาร ACD สำหรับโปรแกรม fconsense	86
3.20 Java class ของบริการท้องถิ่นสำหรับการจัดการอินพุต	87
3.21 Java class สำหรับการอ่านและส่งผลลัพธ์ของรูปแบบ Treefile	88
3.22 Java class สำหรับการอ่านและส่งผลลัพธ์ในรูปแบบ Text	89
3.23 ขั้นตอนการใช้งานบริการท้องถิ่นด้วยโปรแกรม Soaplab Analysis Tool	
3.24 การใช้เครื่องมือ ACD2XML	91
3.25 เวิร์กโฟลว์ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์	92
3.26 เวิร์กโฟลว์ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโน	93
3.27 เวิร์กโฟลว์ใหม่ที่ให้บริการท้องถิ่น (ในกรอบสีดำ)	95
3.28 โครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์ของกุ้ง	98
3.29 โครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโนของกุ้ง	99

รายการรูป (ต่อ)

รูป	หน้า
4.1 Flowchart การทำงานในภาพรวมของเวิร์คโฟลว์สำหรับวิเคราะห์สปีชีส์ของมนุษย์.....	104
4.2 Data flow diagram ของเวิร์คโฟลว์.....	106
4.3 เวิร์คโฟลว์การค้นหาข้อมูลสปีชีส์ที่สนใจ	107
4.4 เวิร์คโฟลว์การค้นหาข้อมูลของยีน	109
4.5 เวิร์คโฟลว์การค้นหาข้อมูล Gene Ontology	110
4.6 เวิร์คโฟลว์การค้นหาข้อมูล OMIM.....	111
4.7 เวิร์คโฟลว์การค้นหาข้อมูล Motif ของยีนจากฐานข้อมูล KEGG	113
4.8 เวิร์คโฟลว์การค้นหา Pathways หรือกลไกการทำงานของยีนจากฐานข้อมูล KEGG	115
4.9 เวิร์คโฟลว์การค้นหาข้อมูลของยีนจากฐานข้อมูล Entrez จาก NCBI	117
4.10 เวิร์คโฟลว์การค้นหาข้อมูลสปีชีส์.....	119
4.11 หน้าต่างแสดงทุกบริการในเวิร์คโฟลว์ทำงานสำเร็จ โดยมีแสดงด้วยบล็อกสีเขียว	122
4.12 หน้าต่างแสดงส่วนการค้น Gene Ontology ทำงานได้สำเร็จแต่ไม่ได้ให้ผลลัพธ์ใดๆ.....	122
5.1 บริการท้องถิ่น ‘XPath From Text’ ในโปรแกรมทาวเวอร์น่า	127
5.2 โครงสร้างของบริการไปโอมาร์ทในสิ่งแวดล้อมของทาวเวอร์น่า.....	128
5.3 ตัวอย่างเอกสาร Scufi ที่มีการบันทึกอิลิเมนต์ของบริการไปโอมาร์ท	129
5.4 Context diagram ของระบบการตรวจสอบของเวิร์คโฟลว์.....	130
5.5 Data flow diagram ชั้นที่ 2 อธิบายกระบวนการตรวจสอบบริการไปโอมาร์ท	131
5.6 Data flow diagram ชั้นที่ 3 อธิบายกระบวนการตรวจสอบบริการไปโอมาร์ท	132
5.7 การใช้งานบริการ XPath_From_Text.....	135
5.8 ตัวอย่างการใช้งานบริการ Get_web_page_from_URL	136
5.9 การกำหนดอินพุตและเอาต์พุตในบีนเชลล์	137

รายการรูป (ต่อ)

รูป	หน้า
5.10 สคริปต์ของบีนเชลล์ในการเปรียบเทียบแอตทริบิวต์ของบริการไปโอมาร์ท	138
5.11 การปรับแต่งกลไกการทำซ้ำของพอร์ตอินพุตในเวิร์คโฟลว์	140
5.12 โครงสร้างของ Scumf ในการปรับแต่งการทำซ้ำของพอร์ตอินพุต	140
5.13 เวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้องของบริการไปโอมาร์ท	142
5.14 เวิร์คโฟลว์อย่างง่ายสำหรับการทดลอง	143
5.15 โครงสร้างภาษา Scumf ของเวิร์คโฟลว์อย่างง่ายสำหรับการทดลอง	145
5.16 รายงานการตรวจสอบบริการไปโอมาร์ทในส่วนของฟิลด์ข้อมูลที่ล้ำสมัย	146
5.17 รายงานการตรวจสอบบริการไปโอมาร์ทในส่วนของฟิลด์ข้อมูลที่ยังทันสมัย	146
6.1 เวิร์คโฟลว์สำหรับการตรวจสอบบริการไปโอมาร์ท	148
6.2 หน้าต่างโปรแกรมพร้อมที่จะตรวจสอบเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์	149
6.3 เวิร์คโฟลว์กำลังตรวจสอบบริการไปโอมาร์ท	150
6.4 ผลการตรวจสอบฟิลด์เตอร์ของบริการไปโอมาร์ทในโปรแกรมทาเวอร์น่า	150
6.5 รายงานการตรวจสอบฟิลด์เตอร์ของบริการไปโอมาร์ท	151
6.6 ผลการตรวจสอบแอตทริบิวต์ของบริการไปโอมาร์ทในโปรแกรมทาเวอร์น่า	151
6.7 รายงานการตรวจสอบแอตทริบิวต์ล้ำสมัยของบริการไปโอมาร์ท	151
6.8 เวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์ที่ได้รับการแก้ไขแล้ว	155
6.9 การตรวจสอบเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์หลังการปรับแก้แล้ว	156
6.10 ตัวอย่างรายงานผลการตรวจสอบของแอตทริบิวต์ซึ่งมีความทันสมัยทุกตัว	156

สัญลักษณ์คำย่อและตัวย่อ

GWAS	Genome Wide Association Study
SNP	Single Nucleotide Polymorphism
SOA	Service Oriented Architecture
XML	Extensible Markup Language
HTTP	Hypertext Transfer Protocol
SOAP	Service Oriented Architecture Protocol
WSDL	Web Services Description Language
UDDI	Universal Description, Discovery and Integration
EBI	European Bioinformatics Institute
SCUFL	Simple Conceptual Unified Flow Language
WABI	Web API for Biology
DDBJ	DNA Data Bank of Japan
BLAST	Basic Local Alignment Search Tool
AME	Advanced Model Explorer
LGPL	GNU Lesser General Public License

บทที่ 1

บทนำ

1.1 ที่มาและความสำคัญของวิทยานิพนธ์

ในอดีตที่ผ่านมานักวิจัยด้านชีวสารสนเทศศาสตร์ ต้องทำงานในลักษณะที่เรียกว่า “คัดลอกแล้ววาง” (Copying and Pasting) กล่าวคือการคัดลอกข้อมูลที่ต้องการนำไปทำงานต่อจากหน้าเว็บหนึ่ง แล้วนำไปวางเพื่อเป็นอินพุตในอีกหน้าเว็บหนึ่ง โดยใช้หลายๆเว็บหรือหลายโปรแกรมเพื่อให้ได้ผลลัพธ์ที่ต้องการ ซึ่งวิธีการแบบนี้ ต้องใช้เวลามาก เสี่ยงต่อความผิดพลาด และไม่สามารถทำงานวิจัยที่ใหญ่ขึ้นได้ เนื่องจากกับความซับซ้อนของงานวิจัยและการไม่มีมาตรฐานร่วมกันของข้อมูล ดังนั้นระบบการทำงานแบบเวิร์คโฟลว์จึงถูกพัฒนาขึ้น เพื่อช่วยในการสร้างเวิร์คโฟลว์ของระบบงาน เป็นการจัดการการเข้าถึงทรัพยากรต่างๆที่กระจายในอินเทอร์เน็ตด้วยเวิร์คโฟลว์และความสามารถในการจัดการแปลงข้อมูลระหว่างมาตรฐานต่างๆ [1][2][3]

โปรแกรมทาเวอร์นา (Taverna) พัฒนาขึ้นในโครงการมายกริด (myGrid) [3][4] ซึ่งเป็นโปรแกรมประยุกต์ที่รวบรวมการเข้าถึงบริการและทรัพยากร ที่กระจายอยู่ในอินเทอร์เน็ตเหล่านั้นเข้าไว้ด้วยกัน เป็นเสมือนหนึ่งว่าอยู่บนระบบเดียวกันด้วยเทคโนโลยีเว็บเซอร์วิส (Web Service) อย่างไรก็ตามงานวิจัยที่ใช้เทคโนโลยีเวิร์คโฟลว์ดังกล่าว ยังจัดว่ามีลักษณะเฉพาะตัวจึงจำเป็นต้องพัฒนาเวิร์คโฟลว์ที่เหมาะสมและพัฒนาเว็บเซอร์วิสเพิ่มเติม เพื่อให้ทำงานได้อย่างถูกต้อง ให้ผลลัพธ์ที่ตรงตามวัตถุประสงค์และมีประสิทธิภาพในการทำงาน

การพัฒนาเวิร์คโฟลว์เชื่อมโยงเว็บเซอร์วิส เพื่อวิเคราะห์และทำนายผลกระทบของสเนป (Single Nucleotide Polymorphism หรือ SNP) ในงานวิจัยเภสัชพันธุศาสตร์ เป็นโครงการที่ร่วมมือกับนักชีวสารสนเทศของหน่วยไวรัสวิทยาและจุลชีววิทยาโมเลกุล ภาควิชาพยาธิวิทยา คณะแพทยศาสตร์ โรงพยาบาลรามาธิบดี มหาวิทยาลัยมหิดล โดยงานวิจัยในลักษณะนี้เรียกว่า Genome Wide Association Study หรือ GWAS [5] เป็นงานวิจัยที่ซับซ้อนมากไม่สามารถทำให้สำเร็จได้อย่างมีประสิทธิภาพและรวดเร็ว ด้วยวิธีการตัดต่อข้อมูลจากอินเทอร์เน็ตโดยวิธีคัดลอกแล้ววางโดยนักวิจัย (Manual) เพียงอย่างเดียวเท่านั้น เนื่องจากข้อมูลสเนป [6] มีจำนวนมากมายถึง

ประมาณ 10 ล้านตำแหน่ง กระจายกันอยู่บนโครโมโซม 23 คู่ของมนุษย์ โดยสปีปจะตอบสนองต่อการใช้ยาและภาวะการที่เป็นโรค [6][7] และข้อมูลที่ใช้ก็กระจายอยู่ที่ต่างๆบนอินเทอร์เน็ต ซึ่งอาจจะมีโอกาสผิดพลาดได้สูงหากให้นักวิจัยตัดต่อข้อมูลด้วยวิธีคัดลอกแล้ววางดังกล่าว ดังนั้นกรณีศึกษาคือการพัฒนาเวิร์คโฟลว์โดยใช้โปรแกรมทาเวอร์นา ในการทำนายผลกระทบอันเนื่องมาจากการเปลี่ยนแปลงของลำดับเบส (SNP) ต่างๆ รวมทั้งสร้างระบบที่มีความสามารถในการคัดเลือกตำแหน่งสปีปที่น่าสนใจในการนำไปศึกษาด้านหน้าที่การทำงาน (Functional Study) หรือผลกระทบต่อโปรตีนของงานด้านชีวสารสนเทศศาสตร์ต่อไป

นอกจากนี้การใช้ฐานข้อมูลจากหลายๆที่ช่วยเลือกยีน หรือจัดลำดับยีนที่มีความน่าสนใจหรือมีความสำคัญเพื่อให้สามารถนำไปศึกษาต่อได้ ซึ่งจะเป็นเรื่องที่ยากมากหากไม่ใช่เวิร์คโฟลว์ เพราะข้อมูลมีขนาดใหญ่และใช้จากหลายๆแหล่งที่กระจายกันอยู่บนอินเทอร์เน็ต การเลือกเว็บเซอร์วิสที่ถูกต้องและตรงตามจุดประสงค์มาใช้ประกอบเป็นเวิร์คโฟลว์นั้น ไม่สามารถจะเลือกได้ทันที เนื่องจากปัจจุบันมีบริการทางด้านชีวสารสนเทศศาสตร์กว่า 30,000 บริการ ซึ่งทำหน้าที่ที่ทั้งที่เหมือนกันและแตกต่างกันไปตามแต่ละบริการนั้นๆ [8][9] ดังนั้นวิทยานิพนธ์นี้จะเกี่ยวข้องกับการวาดเวิร์คโฟลว์และการเลือกเว็บเซอร์วิสที่เหมาะสม โดยให้ข้อเสนอแนะที่เป็นรูปธรรมซึ่งเป็นวัตถุประสงค์แรก

การทดลองสร้าง และการทำงานกับเวิร์คโฟลว์กรณีศึกษาการวิเคราะห์สปีปของมนุษย์ด้านเภสัชพันธุศาสตร์ที่เป็นงานวิจัยร่วมกับมหาวิทยาลัยมหิดล บริการหลักที่เลือกนำมาใช้งานคือบริการไบโอมาร์ท (BioMart) [11][12] และเนื่องจากงานวิจัยมีความซับซ้อนจึงพัฒนาเวิร์คโฟลว์ที่ละส่วนๆตามลักษณะขั้นตอนของการทำงานขั้นตอนนั้นๆ นอกจากนี้เวิร์คโฟลว์หนึ่งๆมีการทำงานในลักษณะ Data flow กล่าวคือ เอาที่พูดของบริการหนึ่งๆ จะเป็นอินพุตให้กับบริการอื่นๆ หากบริการใดบริการหนึ่งผลิตผลลัพธ์ไม่ถูกต้องตรงกับวัตถุประสงค์ของการทำงาน ก็จะทำให้ Data flow ในเวิร์คโฟลว์ส่วนนั้นไม่ถูกต้องตรงตามวัตถุประสงค์ไปด้วย และโปรแกรมทาเวอร์นาเองก็ไม่มีวิธีการดีบักหรือตรวจสอบใดๆเพียงพอที่จะตรวจสอบเวิร์คโฟลว์เพื่อแก้ปัญหานี้ได้จากการทดลองสร้างและทำงานกับเวิร์คโฟลว์ในงานวิจัยที่กล่าวข้างต้น สามารถสรุปสาเหตุที่ทำให้เสียเวลาได้ดังนี้

- 1) หากเวิร์คโฟลว์ใดใช้เวลาในการทำงานที่ยาวนาน ก็ต้องรองกว่าเวิร์คโฟลว์นั้นๆจะทำงานเสร็จสิ้นเสียก่อน จึงจะทราบว่าผลลัพธ์การทำงานนั้นไม่ถูกต้องตามความต้องการ ทำให้เสียเวลามาก เช่น เป็นชั่วโมงๆหรือเป็นวัน เป็นต้น

2) เมื่อทราบว่าผลลัพธ์การทำงานไม่ถูกต้องตามข้อที่ 1 แล้ว ผู้ใช้งานก็จะเริ่มการดีบั๊ก (Debug) หรือตรวจสอบ (Validate) Data flow ของแต่ละบริการเพื่อให้ทราบว่าบริการหรือ โพรเซสเซอร์ (Processor) ใดในเวิร์คโฟลว์ที่ทำงานผิดพลาด

3) เมื่อทราบแล้วว่าบริการ (Processor) ใดทำงานผิดพลาด ผู้ใช้งานจะพยายามตีกรอบการตรวจสอบให้แคบลงโดยการตรวจสอบอย่างถี่ถ้วนที่บริการนั้นๆ ไม่ว่าจะป็นชนิดของข้อมูลหรือไวยากรณ์ของสคริปต์ต่างๆ ซึ่งในบางครั้งก็ถูกต้องดีแล้ว แต่เมื่อทดสอบโดยให้เวิร์ค-โฟลว์ทำงานอีกครั้ง ก็ยังให้ผลผิดพลาดเช่นเดิม

4) เมื่อเป็นดังข้อที่ 1 ถึง 3 ก็พอที่จะสรุปได้ว่า ปัญหาอาจจะไม่ได้เกิดกับ Data flow แต่น่าจะอยู่ที่ตัวบริการเองต่างหาก ดังนั้นผู้ใช้จะเข้าไปตรวจสอบข้อกำหนดการปรับแต่งของบริการไบโอมาร์ทที่เลือกใช้ (Configuration) ซึ่งก็จะพบว่าคุณลักษณะ (ฟิลเตอร์และแอตทริบิวต์) บางค่าของบริการที่เรียกใช้ในเวิร์คโฟลว์นั้น ได้ล้าสมัยหรือไม่มีให้บริการแล้ว

5) วิธีการตรวจสอบของโปรแกรมทาเวอร์นาเอง ก็ไม่เพียงพอที่จะตรวจสอบเวิร์คโฟลว์เพื่อแก้ปัญหาได้ ทำให้ต้องเสียเวลาไปมากเป็นหลายๆชั่วโมงหรือเป็นวัน ยิ่งเวิร์คโฟลว์มีความซับซ้อนทั้งจำนวนบริการไบโอมาร์ทและจำนวนคุณสมบัติมากเท่าไร ก็ยิ่งเสียเวลาเพิ่มขึ้นไปด้วย

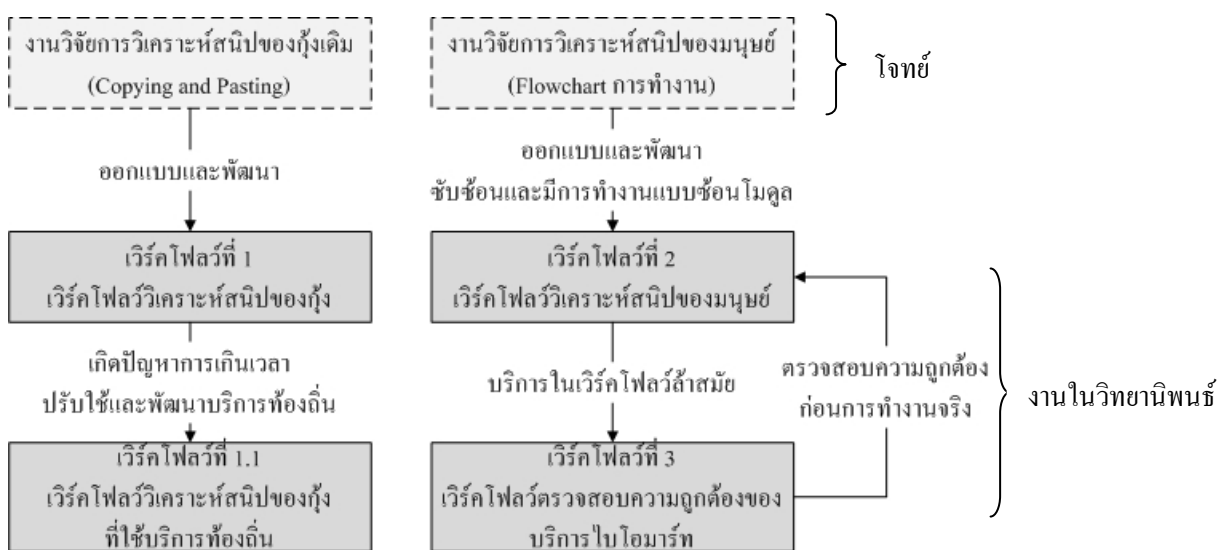
6) สำหรับการแก้ปัญหา Data flow ที่มีความผิดพลาดคือ การวาดเวิร์คโฟลว์ใหม่ โดยการลบบริการที่ล้าสมัยออก และเพิ่มบริการเดิมที่เป็นปัจจุบันลงไป ในเวิร์คโฟลว์ไคอะแกรมที่มีปัญหาดังกล่าว โดยเลือกจากรายชื่อบริการที่ปรากฏในโปรแกรมทาเวอร์นา เพราะการเริ่มต้นของโปรแกรมทาเวอร์นาทุกครั้งนั้น จะมีการค้นหาบริการต่างๆจากข้อมูลที่เป็นค่าโดยปริยายของโปรแกรมเอง (Default) และไบโอมาร์ทก็เป็นส่วนหนึ่งของค่าโดยปริยายนั้นแล้ว ดังนั้นบริการต่างๆที่ปรากฏในรายชื่อบริการ จะเป็นบริการที่ทันสมัยที่สุดที่โปรแกรมนำเสนอต่อผู้ใช้งานในขณะเวลานั้น[8][9]

การค้นหาบริการไบโอมาร์ทใดๆที่มีคุณสมบัติต่างๆที่ล้าสมัย หรือสามารถกล่าวได้ว่าบริการไบโอมาร์ทที่ล้าสมัยนั้น จะต้องทำด้วยวิธีการตรงไปตรงมาโดยผู้ใช้งานเอง (Manual) คือ ผู้ใช้ต้องเรียกหน้าต่างการปรับแต่งของบริการไบโอมาร์ทขึ้นมา แล้วสังเกตว่าคุณสมบัติใดที่เลือกไว้ก่อนหน้านั้น ยังปรากฏให้เลือกได้ในหน้าต่างการปรับแต่งนี้หรือไม่ หากไม่มีแล้ว แสดงว่าคุณสมบัตินั้นๆได้ถูกผู้ให้บริการเอาออกหรือเปลี่ยนไปเป็นชื่อใหม่ เป็นต้น เหล่านี้ล้วนต้องอาศัยประสบการณ์ และการสังเกตจากการใช้บริการไบโอมาร์ทของผู้ใช้งานเอง ยิ่งกว่านั้น จำนวนของบริการไบโอมาร์ทที่อยู่ในเวิร์คโฟลว์, จำนวนคุณสมบัติที่เรียกใช้ในบริการไบโอมาร์ท และความ

ซับซ้อนของเวิร์คโฟลว์ ปัจจัยเหล่านี้ล้วนมีผลต่อเวลาในการตรวจสอบหรือการค้นหาสาเหตุที่ทำให้เวิร์คโฟลว์ทำงานได้ผลลัพธ์ที่ไม่ตรงกับความต้องการ

ดังนั้นวิทยานิพนธ์นี้ จึงเสนอการพัฒนากลไกการตรวจสอบเวิร์คโฟลว์ก่อนการทำงาน โดยจะมุ่งเน้นไปที่การตรวจสอบบริการไปโอมาร์ทในเวิร์คโฟลว์วิเคราะห์สนิป เพราะเป็นบริการหลักที่ใช้ในเวิร์คโฟลว์ โดยเสนอการพัฒนาเวิร์คโฟลว์ภายใต้สิ่งแวดล้อมของโปรแกรมทาวเวอร์น่าที่สามารถตรวจสอบคุณสมบัติของบริการไปโอมาร์ทได้ เพื่อให้ข้อมูลที่มีนัยสำคัญต่อผู้ใช้งาน จะได้มีความสะดวกและง่ายต่อการตรวจสอบ และปรับแก้เวิร์คโฟลว์ซึ่งสามารถลดเวลาที่เสียไปได้อย่างมากเป็นวัตถุประสงค์ที่สอง

อีกประการหนึ่ง ในการสร้างเวิร์คโฟลว์เพื่อวิเคราะห์การเปลี่ยนแปลงลำดับเบส (Single-nucleotide polymorphism หรือ SNP) [6] ในกุ้ง เพื่อประโยชน์ในการพัฒนาสายพันธุ์ซึ่งเป็นความร่วมมือกับสถานวิจัยจีโนมิกส์และชีวสารสนเทศ มหาวิทยาลัยสงขลานครินทร์ พบว่าในกระบวนการวิเคราะห์ข้อมูลลำดับพันธุกรรมของกุ้ง เพื่อสร้างโครงสร้างต้นไม้สายวิวัฒนาการ (Phylogenetic tree) [13][15] มีการวิเคราะห์และถ่ายโอนข้อมูลจำนวนมากในแต่ละบริการ การทำงานมีความล่าช้าจนทำให้มีปัญหาคงการเกินเวลา (Time-out) และความไม่คงเส้นคงวาในการทำงานเกิดขึ้น (Inconsistency) ส่งผลให้เวิร์คโฟลว์ทำงานไม่สำเร็จส่งผลให้ไม่ได้ผลลัพธ์ที่ต้องการ จึงเสนอการสร้างบริการท้องถิ่นในส่วนการทำงานเพื่อสร้างโครงสร้างต้นไม้สายวิวัฒนาการ ในการแก้ไขปัญหาดังกล่าว ซึ่งสามารถช่วยประหยัดเวลาในประมวลผลและให้ผลลัพธ์ที่ถูกต้องตรงตามความต้องการได้ ซึ่งเป็นส่วนหนึ่งของวิทยานิพนธ์นี้ด้วยในวัตถุประสงค์สุดท้าย



รูปที่ 1.1 สรุปงานในวิทยานิพนธ์

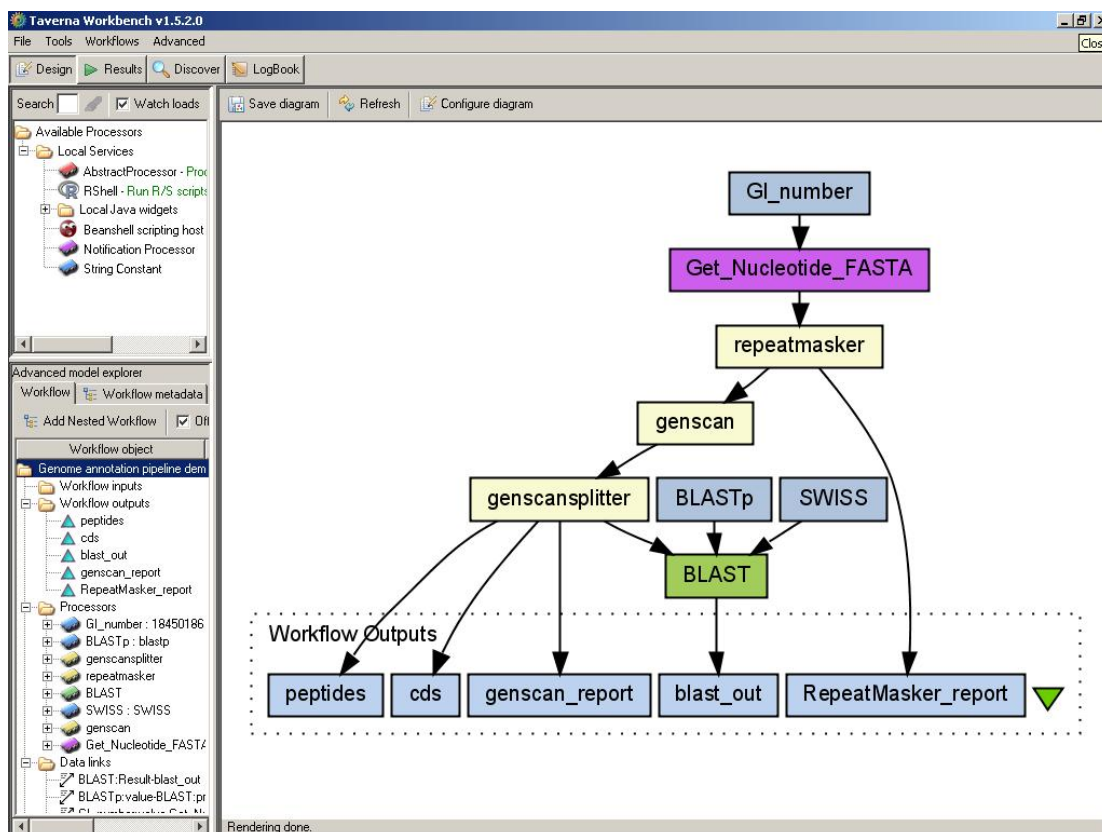
รูปที่ 1.1 แสดงแผนผังสรุปรงานวิทยานิพนธ์ โดยประกอบด้วยการวิเคราะห์ สนิปของกึ่งและของมนุษย์ วิทยานิพนธ์นี้ออกแบบพัฒนาเวิร์คโฟลว์เพื่อวิเคราะห์ของกึ่งได้เป็น เวิร์คโฟลว์ที่ 1 ซึ่งจากการทดสอบการทำงาน พบว่าเวิร์คโฟลว์เกิดปัญหาเกินเวลาและความไม่คง เส้นคงวาในการทำงาน จึงมีการแก้ไขปัญหาดังกล่าว โดยสร้างบริการท้องถิ่นขึ้นมาใช้งานเองใน ส่วนของกระบวนการที่เกิดปัญหาเท่านั้นและพัฒนาได้เป็นเวิร์คโฟลว์ที่ 1.1 ซึ่งสามารถแก้ไข ปัญหาดังกล่าวได้

การวิเคราะห์สนิปของมนุษย์ มีระบบการทำงานความซับซ้อนและมีลักษณะการ งานแบบซ้อนโมดูลๆหรือ Nested workflow วิทยานิพนธ์นี้ออกแบบและพัฒนาเวิร์คโฟลว์เพื่อ วิเคราะห์ของมนุษย์ได้เป็นเวิร์คโฟลว์ที่ 2 แต่พบว่าเวิร์คโฟลว์เกิดปัญหาการล่าช้าของบริการและ ทำงานได้ไม่ตรงตามความต้องการ จึงมีการแก้ไขปัญหาดังกล่าวโดยสร้างระบบตรวจสอบเวิร์ค- โฟลว์ก่อนการทำงานจริงได้เป็นเวิร์คโฟลว์ที่ 3 ซึ่งสามารถแก้ไขปัญหาดังกล่าวได้

1.2 การตรวจเอกสาร (Literatures Review)

1.2.1 โปรแกรมทาวเวอร์นา

การทดลองแบบ *in silico* หรือการทำการทดลองด้วยคอมพิวเตอร์ในงานด้าน วิทยาศาสตร์สุขภาพ (Life Sciences) เป็นการรวบรวมเครื่องมือวิเคราะห์และฐานข้อมูลที่กระจายใน ที่ต่างๆเข้าไว้ด้วยกัน การทดลองแบบนี้ทำได้โดยการ คัดลอกแล้ววาง ระหว่างเว็บเพจหรือโดยการ เขียนโปรแกรมเพื่อปะติดปะต่อทรัพยากรเหล่านี้เข้าไว้ด้วยกันเช่น สคริปต์ของภาษา Perl เป็นต้น วิธีการเหล่านี้มีความจำเป็นมากขึ้นในงานชีวสารสนเทศ กลุ่มงานวิจัยอิสระต่างๆ (Autonomous research) [3][4] ได้สร้างและดูแลฐานข้อมูลที่แตกต่างกันและได้พัฒนาอัลกอริทึมและเครื่องมือ วิเคราะห์ที่มีความสัมพันธ์กัน ทรัพยากรและปริมาณของข้อมูลสาธารณะที่เปิดให้เข้าใช้ได้เพิ่มขึ้น อย่างรวดเร็วในยุคข้อมูลจีโนม ปัจจัยเหล่านี้ทำให้เกิดการรวมตัวกันของกลุ่มวิจัยดังกล่าว อย่างไรก็ตามยังมีปัญหาเรื่องมาตรฐานการใช้ข้อมูลร่วมกัน ความแตกต่างกันของระบบที่มีอยู่เดิม และ ลักษณะทางด้านภูมิศาสตร์ที่ตั้งของบริการข้อมูลในอินเทอร์เน็ต



รูปที่ 1.2 โปรแกรมทาวเวอร์นาแสดงเวิร์คโฟลว์และเว็บเซอร์วิสโดยการแบ่งประเภทด้วยสีต่างๆ [4]

โปรแกรมทาวเวอร์นา (Taverna) (รูปที่ 1.2) พัฒนาขึ้นในโครงการมายกริดซึ่งเป็นโปรแกรมประยุกต์ที่รวบรวมการเข้าถึงบริการและทรัพยากร ที่กระจายอยู่ในอินเทอร์เน็ตเข้าไว้เป็นเสมือนหนึ่งว่าอยู่บนระบบเดียวกันด้วยเทคโนโลยีเว็บเซอร์วิส อย่างไรก็ตามงานวิจัยที่ใช้เทคโนโลยีเวิร์คโฟลว์ดังกล่าว ยังมีลักษณะเฉพาะตัวที่ไม่สามารถใช้เวิร์คโฟลว์ตามบริการตัวอย่างเวิร์คโฟลว์ที่มีให้ จึงจำเป็นต้องพัฒนาเวิร์คโฟลว์ที่เหมาะสม และพัฒนาเว็บเซอร์วิสเพิ่มเติม เพื่อให้ทำงานได้อย่างถูกต้อง ให้ผลลัพธ์ที่ตรงตามวัตถุประสงค์และมีประสิทธิภาพในการทำงาน [2][10]

1.2.2 โปรแกรมทาวเวอร์นาสอง รุ่นทดลองใช้

โปรแกรมทาวเวอร์นาตั้งแต่รุ่น 1.7.0 มีปลั๊กอินชื่อว่า *Taverna 2 preview* ซึ่งเป็นการทดลองโปรแกรมทาวเวอร์นา รุ่น 2 เพื่อการทำงานที่เรียกว่า *Workflow health check* หรือการตรวจสอบเสถียรภาพของเวิร์คโฟลว์ โดยการตรวจสอบว่า เอ็นพอยต์ของเว็บเซอร์วิสต่างๆในเวิร์ค

โพล์มีตอบกลับผ่าน HTTP (Responding) หรือไม่ [8][9] หากมีการตอบกลับก็คือสามารถเข้าถึงได้ หากไม่มีการตอบกลับ ก็ไม่สามารถที่จะเปิดเวิร์คโพล์ขึ้นมาทำงานได้ ปัญหาที่อาจจะมาจากเครือข่ายของผู้ใช้งานเองไม่สามารถติดต่อภายนอกได้ หรือไม่ก็มีปัญหาที่ฝั่งของผู้ให้บริการเองอันเนื่องมาจากหลายๆสาเหตุ

อีกประการหนึ่งในการออกแบบและการสร้างเวิร์คโพล์ ในงานวิจัยด้านเภสัช-พันธุศาสตร์ซึ่งเป็นงานวิจัยที่มีขนาดใหญ่และมีกระบวนการที่ซับซ้อน จำเป็นต้องออกแบบการทำงานเป็นโมดูลๆทำงานในกระบวนการที่เกี่ยวข้องซึ่งมีการเรียกใช้โมดูลอื่นๆด้วย ในสิ่งแวดล้อมของทาวเวอร์นี้เรียกโมดูลการทำงานของแต่ละกระบวนการเหล่านี้ว่า Nested workflow และพบว่า *Taverna 2 preview* ไม่สามารถตรวจสอบ Nested workflow ได้

ดังนั้น *Taverna 2 preview* ไม่สามารถตรวจจับบริการที่ล่าสมัยไปแล้วและไม่สามารถตรวจสอบความเสถียรภาพของเวิร์คโพล์ที่มีความซับซ้อนได้ การแก้ไขข้อจำกัดของ *Taverna 2 preview* เหล่านี้ เป็นส่วนหนึ่งของวิทยานิพนธ์นี้ด้วย

1.2.3 บริการไบโอมาร์ท (BioMart service)

ไบโอมาร์ท (BioMart) [11][12] คือระบบการจัดการข้อมูลเชิงคิวรีที่พัฒนาโดยสถาบันวิจัยมะเร็ง Ontario (OICR) และสถาบันชีวสารสนเทศแห่งยุโรป (EBI) ระบบไบโอมาร์ทมีข้อมูลหลากหลายชนิดที่นำเสนอต่อผู้ใช้งาน เหมือนระบบเหมืองข้อมูลซึ่งทำหน้าที่ในการค้นหาข้อมูลได้อย่างซับซ้อน โดยมีหน้าต่างการทำงานเหมือนหน้าเว็บไซต์ที่เรียกว่า ‘Out of the box’ ที่สามารถติดตั้งและปรับแต่งบริการได้ตามความต้องการของผู้ใช้งาน กลไกเหล่านี้มีลักษณะแบบกราฟิกส์และข้อความ โดยมีพื้นฐานการทำงานแบบโปรแกรมประยุกต์ หรือใช้เทคโนโลยีเว็บเซอร์วิสหรือ Application Programming Interface (API) ที่เขียนด้วยภาษา Perl และ Java

ไบโอมาร์ทสร้างขึ้นเพื่อสนับสนุนการคิวรีข้อมูลได้เป็นอย่างดี และนอกจากนี้ยังสามารถปรับแต่งให้ทำงานกับเซิร์ฟเวอร์ฐานข้อมูล Distributed Annotation System (DAS) 1.5 ซึ่งใช้ในการแลกเปลี่ยนคำจำกัดความประกอบของลำดับจีโนมหรือโปรตีนได้ด้วย [16][17] โดยกระบวนการการแปลงจากแหล่งข้อมูลไปเป็นรูปแบบของไบโอมาร์ทนั้น สามารถทำได้อย่างอัตโนมัติด้วยเครื่องมือที่ไบโอมาร์ทได้เตรียมไว้ให้แล้ว ซึ่งสนับสนุนระบบจัดการฐานข้อมูลเชิงความสัมพันธ์ทั้ง MySQL [18], Oracle [19] และ Postgres [20] ไบโอมาร์ทเป็นซอฟต์แวร์แบบเปิด

ภายใต้ลิขสิทธิ์ของ GNU Lesser General Public License (LGPL) [21] ซึ่งผู้ใช้สามารถเข้าถึงและใช้งานโดยไม่มีข้อผูกมัดหรือเงื่อนไขใดๆ

บริการไบโอมาร์ทหนึ่งๆจะประกอบด้วยชื่อของบริการ (Service Name) และเนื่องจากบริการไบโอมาร์ททำงานในลักษณะเหมือนข้อมูล ดังนั้นบริการของไบโอมาร์ทสามารถคิวรีข้อมูลจากดาตาเซต (Dataset) ที่ผู้ใช้บริการกำหนดไว้ได้ โดยในแต่ละดาตาเซตก็จะประกอบด้วยคุณสมบัติต่างๆที่แตกต่างกันออกไปคือ ฟิลเตอร์ หรือฟิลด์ข้อมูลสำหรับเป็นเงื่อนไขการคิวรี (Filter or query filed) และแอตทริบิวต์ หรือฟิลด์ข้อมูลสำหรับเป็นผลลัพธ์ของการคิวรี (Attribute or output filed) [22] ในการออกแบบเวิร์คโฟลว์เพื่อวิเคราะห์สnipของมนุษย์พบว่าบริการไบโอมาร์ทหนึ่งๆสามารถให้บริการข้อมูลจากฐานข้อมูลเพียงหนึ่งฐานข้อมูล อย่างพอเพียงต่อความต้องการ และการทำงานแบบหนึ่งบริการต่อหนึ่งฐานข้อมูล ยังสะดวกต่อการพัฒนาและการปรับปรุงเวิร์คโฟลว์ได้ในอนาคต ดังนั้นงานวิจัยนี้จึงได้ตีกรอบเนื้อหาของบริการไบโอมาร์ทไว้ว่า บริการหนึ่งๆจะมีหนึ่งชื่อบริการและสามารถคิวรีข้อมูลได้หนึ่งฐานข้อมูล ซึ่งฐานข้อมูลหนึ่งๆประกอบด้วยหลายๆคุณสมบัติ ดังแสดงโครงสร้างบริการไบโอมาร์ทในรูปที่ 1.3

A BioMart Service	
A Processor/Service Name	
A Dataset	
Attributes	Filters

รูปที่ 1.3 โครงสร้างบริการไบโอมาร์ทในงานวิจัย

บริการไบโอมาร์ทเป็นปลั๊กอินตัวหนึ่งในโปรแกรมทาเวอร์น่า ทำให้สามารถใช้บริการข้อมูลได้อย่างง่ายด้วยเทคโนโลยีเว็บเซอร์วิส การเข้าถึงข้อมูลต่างๆจะเป็นลักษณะการสื่อสาร Java Database Connectivity (JDBC) [23] ซึ่งเป็น API ใน Java ที่ใช้สำหรับติดต่อไปยังฐานข้อมูลกลางของไบโอมาร์ทและนำเสนอต่อผู้ใช้งานในลักษณะของกราฟิกส์ สามารถติดตั้งและปรับแต่งบริการได้ตามความต้องการของผู้ใช้งาน [2]

1.2.4 การพัฒนาเวิร์คโฟลว์สำหรับการวิเคราะห์สลับของมนุษย์

หน่วยไวรัสวิทยาและจุลชีววิทยาโมเลกุล คณะแพทยศาสตร์ โรงพยาบาลรามาธิบดี มหาวิทยาลัยมหิดล ศึกษาสลับที่เกี่ยวข้องกับการแพ้ยาเนวีราพีน (Nevirapine) [24] ซึ่งเป็นยาด้านเชื้อไวรัสเอชไอวี โดยเป็นยาในกลุ่ม Non-nucleoside reverse transcriptase inhibitor (NNRTI) [24] การแพ้ยาชนิดนี้จะมีอาการตั้งแต่เป็นผื่นคันเล็กน้อยจนถึงอาการหนัก และอาจทำให้เกิดโรค Steven's Johnson Syndrome [25] ซึ่งจะทำให้เกิดความลำบากในการใช้ชีวิตกับตัวผู้ติดเชื้อ และยังทำให้เสียเงินจำนวนมากมายในการรักษาอาการแพ้ยา ในท้ายที่สุดการแพ้ยาอาจทำให้ถึงแก่ชีวิตได้ การตรวจการเปลี่ยนแปลงของลำดับเบสหรือสลับ (SNP) ที่มีความเกี่ยวข้องหรือสัมพันธ์กับกลไกการแพ้ยา จะเป็นงานบริการที่สำคัญในอนาคตกับผู้ติดเชื้อเอชไอวีที่ต้องรับประทานยา Nevirapine โดยการศึกษาสลับที่เกี่ยวข้องกับการแพ้ยา Nevirapine นี้จะใช้วิธีการทำแบบ GWAS มีกลุ่มการศึกษาแบ่งเป็น 2 กลุ่ม คือ กลุ่มที่ติดเชื้อเอชไอวีและได้รับประทานยาเนวีราพีนแต่ไม่เกิดอาการแพ้ยา และกลุ่มที่ติดเชื้อเอชไอวีและได้รับประทานยาเนวีราพีนและมีอาการแพ้ยาเกิดขึ้น ขณะนี้การศึกษาอยู่ในขั้นตอนของการวิเคราะห์ผลทางสถิติจากห้องปฏิบัติการ คาดว่าหลังจากที่วิเคราะห์ผลทางสถิติจากห้องปฏิบัติการเสร็จแล้ว จะยังคงมีตำแหน่งสลับที่ต้องทำการศึกษาลงลึกไปในรายละเอียดต่อไปเพื่อหาตำแหน่งสลับที่แท้จริง รวมทั้งกลไกที่เกี่ยวข้องกับการแพ้ยาเนวีราพีนนี้

Lude Franke และคณะ [26] ได้สร้างโปรแกรมคอมพิวเตอร์ที่มีชื่อว่า โปรแกรม Prioritizer โดยเป็นโปรแกรมที่ใช้ในการวิเคราะห์ผล Positional candidate gene prioritization โดยโปรแกรมจะทำหน้าที่นำตำแหน่งเริ่มต้นและตำแหน่งสิ้นสุดของยีน หรือส่วนหนึ่งส่วนใดบนโครโมโซมที่สนใจ ไปเปรียบเทียบกับฐานข้อมูลต่างๆที่มีอยู่ในตัวโปรแกรม ประกอบด้วย ฐานข้อมูล Gene ontology, Protein-Protein interaction และ Microarray co-expression และแสดงผลออกมาให้ผู้ใช้งานเห็นถึงความสัมพันธ์ที่เกี่ยวข้องกับข้อมูลที่ได้ใส่เข้าไปในระบบ วิธีการนี้จะทำให้สามารถคัดกรองจำนวนยีนหรือกลุ่มของยีนที่น่าจะนำไปศึกษาออกเป็นลำดับแรกได้ ในปัจจุบันทางผู้สร้างโปรแกรมได้นำเอาข้อมูลทางด้านสลับและ Genetic linkage จากการสแกนจีโนมทั้งหมด (Whole genome scan) มาประกอบเข้าด้วยกันกับตัวโปรแกรมเดิม เพื่อเพิ่มขีดความสามารถในการคัดกรองตำแหน่งของสลับ หรือยีนที่เกี่ยวข้องกับการเกิดโรคหรือการแพ้ยาชนิดต่างๆ

Young C. Song และคณะ [7] ได้สร้างเวิร์คโฟลว์สำหรับข้อมูล Biological pathway (กลไกการทำงานระหว่างยีน) และ Non-synonymous SNPs ซึ่งใช้ข้อมูลจากเว็บเซอร์วิส BioMoby และใช้โปรแกรมทาเวอร์น่าในการสร้างและออกแบบให้ได้มาเป็นเวิร์คโฟลว์ โปรแกรมที่ Young C. Song และคณะทำการสร้างขึ้นมีชื่อว่า โปรแกรม DataBiNS [7] หน้าที่หลักของโปรแกรมคือรวบรวมข้อมูล Non-synonymous SNPs จากฐานข้อมูล dbSNP [27] กับข้อมูลที่เกี่ยวข้องกับกลไกการทำงานระหว่างยีน (Biological pathway) ต่างๆจากฐานข้อมูล KEGG [28] จากนั้นจะดึงข้อมูลรายการผลงานที่ถูกตีพิมพ์, Gene ontology, Gene annotation และข้อมูล nsSNP [29] ของแต่ละยีนใน Biological pathway ต่างๆนั้น นอกจากนี้โปรแกรมนำเอาโปรแกรม LS-SNP [7][30] ที่มีความสามารถในการทำนายผลกระทบของการเปลี่ยนแปลงสลับที่ที่อาจส่งผลถึงการทำงานของโปรตีนมารวมไว้ด้วยกัน ข้อมูลทั้งหมดจะถูกดึงมาแสดงภายในการสั่งงานเพียงครั้งเดียว สิ่งที่ได้จากโปรแกรมคือ โปรแกรมจะช่วยลดหรือคัดกรองจำนวนและรวมทั้งตำแหน่งของสลับที่มีความเกี่ยวข้องกับกลไกการทำงานระหว่างยีนน้อยๆออก โดยเหลือเฉพาะตำแหน่งที่สามารถนำไปวิเคราะห์ต่อในระดับหน้าที่การทำงานได้

การเลือกตำแหน่งสลับที่น่าสนใจจากผลที่ได้จากสลับชิป รวมทั้งการหาความสัมพันธ์หรือการทำนายถึงผลกระทบ อันเนื่องมาจากการเปลี่ยนแปลงของสลับในมนุษย์กับการตอบสนองต่อยาหรือภาวะการเป็นโรค เป็นสิ่งที่จำเป็นต้องทำหลังจากที่ได้ผลการวิเคราะห์ทางห้องปฏิบัติการมาแล้ว ข้อมูลที่มีมากมายในปัจจุบันจากหลายฐานข้อมูลที่เปิดให้ใช้งานได้ขณะนี้สามารถทำให้เกิดองค์ความรู้ใหม่ๆในด้านการศึกษาสลับ อย่างไรก็ตามผลจากการค้นหาข้อมูลจากฐานข้อมูลในแต่ละฐานข้อมูลดังกล่าวไม่ได้อยู่ในรูปแบบที่สามารถใช้งานได้ง่าย หรือฐานข้อมูลเดียวอาจไม่เพียงพอในการสร้างเป็นองค์ความรู้ได้ คณะผู้วิจัยจึงมีวัตถุประสงค์ที่จะพัฒนาวิธีการทำนายผลกระทบอันเนื่องมาจากการสลับต่างๆ รวมทั้งสร้างระบบที่มีความสามารถในการคัดเลือกตำแหน่งสลับที่น่าสนใจในการนำไปศึกษาต่อทางด้านหน้าที่การทำงาน (Functional Study) ต่อไป ซึ่งยังไม่มีผู้ใดเคยทำมาก่อน คณะผู้วิจัยจะทำการศึกษาสลับที่ละเอียดครอบคลุมและใช้งานง่ายโดยการแก้ปัญหาทางงานวิจัยนี้จะใช้โปรแกรมทาเวอร์น่า ที่มีความสามารถในการสร้างและจัดการรูปแบบข้อมูล รวมทั้งมีความสามารถในการประสานการทำงานกับเว็บเซอร์วิสต่างๆ และแสดงผลออกมาในรูปแบบที่ง่ายต่อการนำไปศึกษาหรือใช้งานต่อไป จึงเป็นเหตุให้เกิดความร่วมมือในการวิจัยระหว่างมหาวิทยาลัยและเป็นกรณีศึกษาของวิทยานิพนธ์นี้

1.3 วัตถุประสงค์

1.3.1. พัฒนาเวิร์คโฟลว์และสร้างการเชื่อมต่อในเวิร์คโฟลว์เพื่อแก้ไขปัญหาในงานวิจัยการทำนายผลกระทบของสปีปในงานวิจัยทางด้านเกษตรพันธุศาสตร์ และการวิเคราะห์สปีปของกึ่ง โดยใช้เทคโนโลยีมายกริด, โปรแกรมทาเวอร์น่า, เว็บเซอร์วิสและใช้เวิร์คโฟลว์ทั้งสองนี้เป็นกรณีศึกษาของงานวิทยานิพนธ์

1.3.2. นำเสนอกลไกในการลดเวลาในการพัฒนา และทดสอบเวิร์คโฟลว์โดยพัฒนาการตรวจสอบความถูกต้องของเวิร์คโฟลว์ก่อนการทำงานจริง ซึ่งจะมุ่งเน้นไปที่การตรวจสอบบริการไบโอมาร์ทภายใต้สิ่งแวดล้อมของโปรแกรมทาเวอร์น่า ที่สามารถตรวจสอบคุณสมบัติของบริการไบโอมาร์ทได้ว่าทันสมัยพร้อมทำงานหรือล้าสมัยไปแล้ว ซึ่งสามารถระบุบริการที่ล้าสมัยไปแล้วในเวิร์คโฟลว์ และให้ข้อมูลที่มีนัยสำคัญต่อผู้ใช้งานจะได้มีความสะดวกและง่ายต่อการปรับแก้เวิร์คโฟลว์ซึ่งสามารถลดเวลาที่เสียไปได้อย่างมาก

1.3.3. นำเสนอกลไกในการลดปัญหาการเกินเวลา และความไม่คงเส้นคงวาในการทำงานที่ส่งผลให้เวิร์คโฟลว์ทำงานไม่สำเร็จโดยพัฒนาบริการท้องถิ่นขึ้นมาใช้งานเอง (Local web services)

1.4 ขอบเขตของงานวิจัย

1.4.1. พัฒนาเวิร์คโฟลว์โดยใช้โปรแกรมทาเวอร์น่ารุ่น 1.5.0 – 1.7.1 และมีกรณีศึกษาคือเวิร์คโฟลว์ในการวิเคราะห์สปีปของกึ่งและของมนุษย์

1.4.2. พัฒนาการตรวจสอบเวิร์คโฟลว์ก่อนการทำงาน โดยจะมุ่งเน้นไปที่การตรวจสอบบริการไบโอมาร์ทภายใต้สิ่งแวดล้อมของโปรแกรมทาเวอร์น่า โดยสามารถตรวจสอบคุณสมบัติของบริการไบโอมาร์ทได้ โดยมีข้อกำหนดในการตรวจสอบเป็น 1 บริการใช้ 1 ฐานข้อมูล

1.4.3. เพิ่มบริการท้องถิ่น (Local web services) เพื่อลดเวลาในการทำงานของเวิร์คโฟลว์ด้วยเทคโนโลยีมายกริดและเว็บเซอร์วิส

1.5 ประโยชน์ที่คาดว่าจะได้รับ

1.6.1. ได้วิธีการแก้ไขปัญหาลดความล่าช้าในการทำงานของเว็บเซอร์วิส ด้านชีวสารสนเทศโดยการใช้บริการท้องถิ่นบนโปรแกรมทาเวอร์น่า

1.6.2. ได้เวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท ภายใต้สิ่งแวดล้อมของทาเวอร์น่าโดยสามารถตรวจสอบคุณสมบัติของบริการไบโอมาร์ทได้

1.6.3. ได้เวิร์คโฟลว์ในการแก้ปัญหาในงานวิจัย การทำนายผลกระทบของการเปลี่ยนแปลงลำดับเบส (SNP) ของมนุษย์ในงานวิจัยเภสัชพันธุศาสตร์

1.6.4. ได้องค์ความรู้เกี่ยวกับการพัฒนาเวิร์คโฟลว์ ด้านชีวสารสนเทศภายใต้สิ่งแวดล้อมของทาเวอร์น่า

1.6 สรุป

วิทยานิพนธ์นี้เป็นการประยุกต์ใช้เทคโนโลยีมายกริด และเว็บเซอร์วิสเพื่อแก้ไขปัญหาในงานวิจัยทางด้านชีวสารสนเทศศาสตร์ ซึ่งเป็นงานวิจัยที่ทำร่วมกับหน่วยงานวิจัยที่มีความชำนาญทางด้านชีวสารสนเทศศาสตร์ของมหาวิทยาลัยสงขลานครินทร์และ โรงพยาบาลรามารชิบตี คณะแพทยศาสตร์ มหาวิทยาลัยมหิดล

สำหรับบทต่อไปจะกล่าวถึงภาพรวม และสถาปัตยกรรมของมายกริดในงานวิทยานิพนธ์ตลอดจนการอธิบายรายละเอียดของบริการไบโอมาร์ท บทที่ 3 จะกล่าวถึงการเลือกบริการที่เหมาะสม โดยใช้เวิร์คโฟลว์การสร้างโครงสร้างต้นไม้สายวิวัฒนาการของกุ้ง ซึ่งเป็นการทำวิจัยร่วมกับศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์ เป็นกรณีศึกษาแรก

บทที่ 4 จะกล่าวถึงการออกแบบและพัฒนาเวิร์คโฟลว์สำหรับวิเคราะห์สปีปของมนุษย์ ซึ่งเป็นการทำวิจัยร่วมกับศูนย์ความเป็นเลิศด้านชีววิทยาศาสตร์ของประเทศไทยหรือ Thailand Center of Excellence for Life Sciences (TCELS) โรงพยาบาลรามาธิบดี มหาวิทยาลัยมหิดล การออกแบบและพัฒนาเวิร์คโฟลว์สำหรับตรวจสอบความถูกต้องของเวิร์คโฟลว์นี้จะกล่าวในบทที่ 5 ผลการทดลองจะกล่าวในบทที่ 6 และสรุปงานวิทยานิพนธ์ในบทที่ 7

บทที่ 2

สถาปัตยกรรม โปรแกรมทาเวอร์นาและเวิร์คโฟลว์

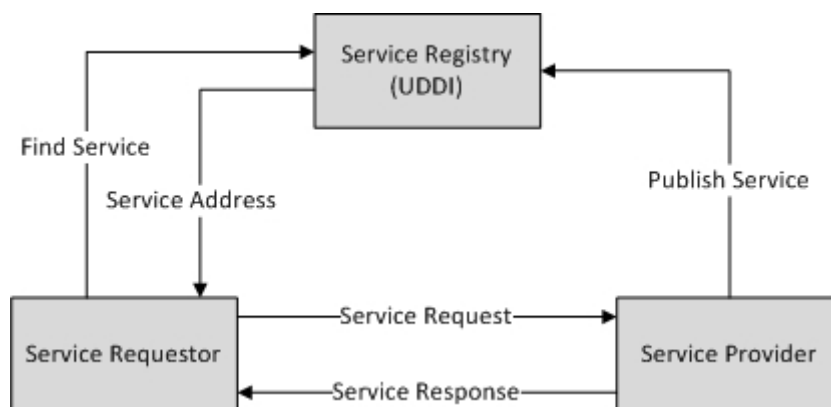
บทนี้จะกล่าวถึงคือสถาปัตยกรรมของงานวิจัย และหลักการสำคัญที่เกี่ยวข้องกับเวิร์คโฟลว์, โปรแกรมทาเวอร์นา, เว็บเซอร์วิสและบริการต่างๆที่จะนำมาสร้างเป็นเวิร์คโฟลว์รวมทั้งรายละเอียดเกี่ยวกับบริการไอโบบมาร์ท

2.1 สถาปัตยกรรมที่เกี่ยวข้องกับวิทยานิพนธ์

งานวิจัยในวิทยานิพนธ์นี้ได้ประยุกต์ใช้สถาปัตยกรรมต่างๆดังนี้

2.1.1 การพัฒนาระบบบนเครือข่ายกริดด้วยแนวความคิดของสถาปัตยกรรมเชิงบริการ

การพัฒนาระบบซอฟต์แวร์เพื่อใช้งานกันระหว่างองค์กรบนเครือข่ายกริด มีปัจจัยที่ต้องนำมาใช้ประกอบการพิจารณาอยู่มากมายด้วยกัน และปัจจัยหลักที่ต้องพิจารณาเป็นอันดับต้นๆก็คือ วิธีการออกแบบระบบที่สามารถเข้ากันได้กับเทคโนโลยีของทุกๆองค์กรที่ร่วมมือกัน รวมไปถึงต้องสามารถเข้ากันได้กับวัฒนธรรมขององค์กรนั้นๆด้วย โดยคำว่าวัฒนธรรมขององค์กรนั้นเกี่ยวข้องกับวิธีการและขั้นตอนในการดำเนินงานของแต่ละองค์กรมีความแตกต่างกัน และระบบที่ใช้อยู่เดิมของแต่ละองค์กรนั้น (Legacy system) มีความแตกต่างกันโดยสิ้นเชิงเป็นต้น ซึ่งปัจจัยทางด้านวัฒนธรรมนี้เอง ได้ผลักดันให้เกิดรูปแบบการพัฒนาระบบที่ตั้งอยู่บนฐานของสถาปัตยกรรมเชิงบริการหรือ Service Oriented Architecture (SOA) [31] โดยองค์ประกอบของ SOA มี 3 ส่วนคือผู้ใช้บริการ (Service Requestor) ผู้ให้บริการ (Service Provider) และฝ่ายทะเบียนของบริการ (Service Registry) [31] โดยผู้ให้บริการจะเป็นผู้จัดสรรบริการให้แก่ผู้ใช้บริการ และฝ่ายทะเบียนของบริการเปรียบเสมือนสมุดหน้าเหลืองสำหรับค้นหาบริการที่จัดสรรโดยผู้ให้บริการ ซึ่งผู้ใช้บริการสามารถค้นหาบริการที่ต้องการได้จากหน่วยนี้ดังรูปที่ 2.1



รูปที่ 2.1 การทำงานของสถาปัตยกรรมเชิงบริการ [30]

สถาปัตยกรรมเชิงบริการ ประกอบด้วยกลุ่มของบริการที่ติดต่อสื่อสารระหว่างกันได้ การติดต่อสื่อสารนี้สามารถส่งผ่านข้อมูลระหว่างบริการสองบริการหรือมากกว่า เพื่อประสานงานกันในการทำงานต่างๆ เริ่มแรกนั้นสถาปัตยกรรมเชิงบริการคือการใช้ Distributed Component Object Model (DCOM) หรือ Object Request Broker (ORBs) ที่อยู่บนพื้นฐานข้อกำหนดของ Common Object Request Broker Architecture (CORBA) [32]

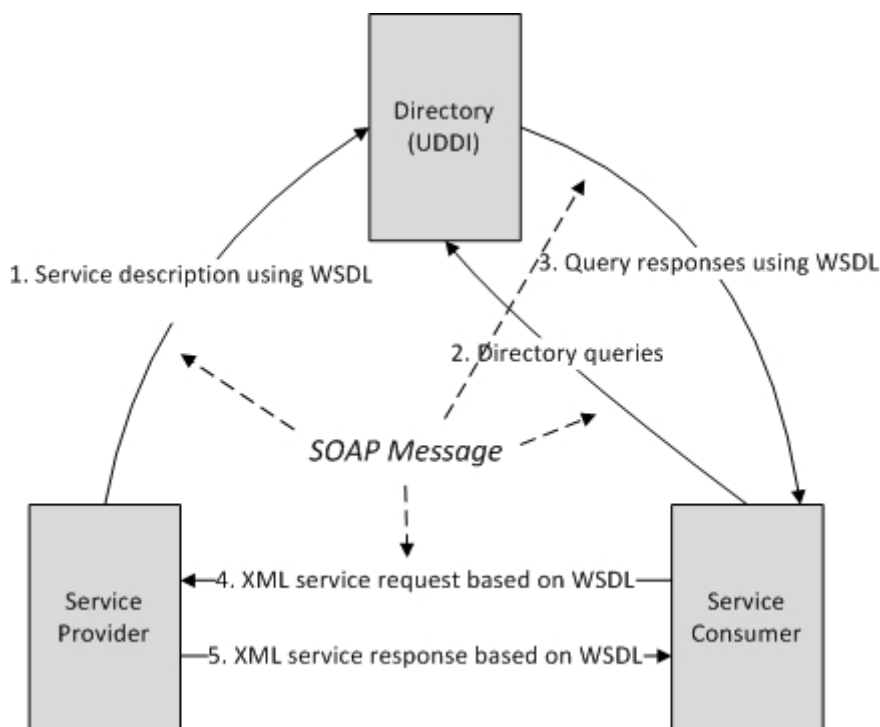
2.1.2 เว็บเซอร์วิส

เว็บเซอร์วิส (Web Service) [32] คือโปรแกรมประยุกต์ที่สร้างให้รองรับการเรียกใช้งานจากโปรแกรมประยุกต์อื่นบนอินเทอร์เน็ต โดยสื่อสารกันด้วยข้อมูลที่อยู่ในรูปแบบ Extensible Markup Language (XML) รูปแบบ XML ที่ใช้นี้กำหนดเป็นมาตรฐานชื่อว่า Simple Object Access Protocol (SOAP) [33] โดยข้อมูลจะถูกส่งผ่านทางโปรโตคอล Hypertext Transfer Protocol (HTTP), Simple Mail Transfer Protocol (SMTP) หรือ File Transfer Protocol (FTP) แต่ที่นิยมใช้มาก คือ HTTP เว็บเซอร์วิสประกอบด้วยดังนี้

- แอ็พพลิเคชัน คือ โปรแกรมประยุกต์ที่ทำหน้าที่ให้บริการอยู่บนเครื่องแม่ข่ายที่เปิดให้บริการตลอดเวลา สามารถติดต่อด้วยโปรโตคอล HTTP และพัฒนาด้วยภาษาที่มีความสามารถในการจัดการกับ SOAP ได้
- SOAP คือ โปรโตคอลหรือระเบียบวิธีในการติดต่อสื่อสารระหว่างเว็บเซอร์วิส โดยใช้ข้อมูลที่กำหนดรูปแบบด้วยภาษา XML ทำให้เว็บเซอร์วิส

สามารถสื่อสารกันได้ แม้ว่าจะอยู่บนเครื่องคอมพิวเตอร์คนละแพลตฟอร์ม หรือพัฒนาด้วยภาษาโปรแกรมที่แตกต่างกัน เมื่อแอปพลิเคชันต้องการใช้งานเว็บเซอร์วิส ต้องเขียนโปรแกรมเพื่อติดต่อกับ SOAP ในภาษาที่ตนใช้ จากนั้น SOAP ก็สร้าง SOAP message เพื่อติดต่อกับแอปพลิเคชันปลายทางให้โดยอัตโนมัติ

- Web Services Description Language (WSDL) [34] คือ เอกสาร XML ที่อธิบายรายละเอียดในการติดต่อกับเว็บเซอร์วิส เพื่อให้แอปพลิเคชันที่ต้องการเรียกใช้เว็บเซอร์วิสรู้ว่า เซอร์วิสนั้นให้บริการอะไรบ้าง และจะติดต่อดีอย่างไร
- UDDI (Universal Description, Discovery and Integration) คือ ไคลเรทอรีที่เก็บรวบรวมเว็บเซอร์วิสที่มีการลงทะเบียนไว้ UDDI จะเก็บรวบรวมข้อมูลของเว็บเซอร์วิสต่างๆไว้ในรูปแบบ WSDL หน้าของ UDDI จะคล้ายกับเว็บไคลเรทอรีหรือเว็บสารบัญ กล่าวคือ UDDI ช่วยให้ผู้พัฒนาเว็บเซอร์วิสได้ประกาศหรือประชาสัมพันธ์บริการของตนเองสู่สาธารณะ และช่วยให้ผู้ใช้งานเว็บเซอร์วิสค้นพบเว็บเซอร์วิสที่ต้องการใช้งานได้
- รูปที่ 2.2 อธิบายการทำงานของเว็บเซอร์วิสได้ดังนี้
 - 1) ผู้ให้บริการจะให้คำอธิบายการใช้บริการโดยใช้ WSDL ไว้ที่ไคลเรทอรีของเว็บเซอร์วิสหรือ UDDI
 - 2) ผู้ขอใช้บริการจะค้นหาที่ต้องการจาก UDDI ทำให้ทราบว่าบริการอยู่ที่ไหน และจะติดต่อกับบริการนั้นได้อย่างไร
 - 3) UDDI จะตอบกลับด้วยเอกสาร WSDL กลับมาให้ผู้ขอใช้บริการโดยจะบอกว่าบริการที่จะเรียกใช้ต้องร้องขอ และจะได้ผลลัพธ์อย่างไร
 - 4) ผู้ขอใช้บริการใช้ WSDL ส่งการร้องขอการใช้งานไปยังผู้ให้บริการ
 - 5) ผู้ให้บริการตอบกลับตามที่ผู้ขอใช้บริการร้องขอ



รูปที่ 2.2 สถาปัตยกรรมการทำงานของเว็บเซอร์วิส [32]

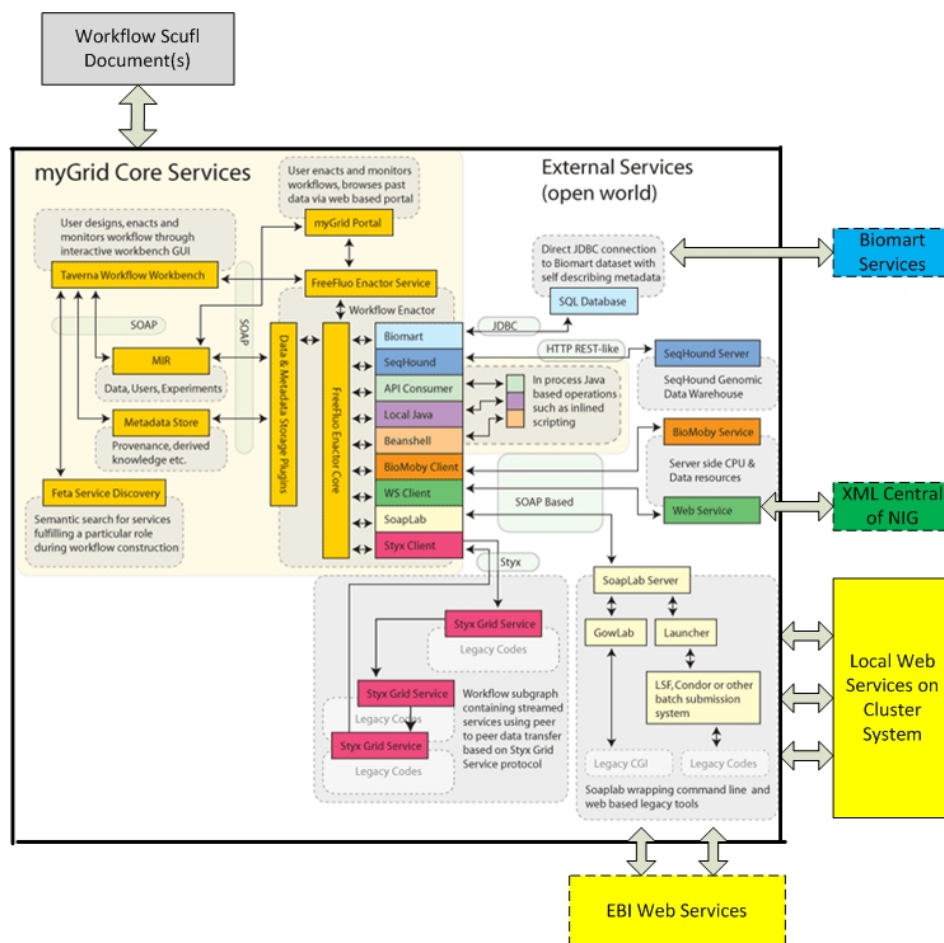
2.1.3 สถาปัตยกรรมของโปรแกรมทาเวอร์นา

รูปที่ 2.3 เป็นสรุปการติดต่อสื่อสารและเชื่อมโยงของคอมพิวเตอร์อันหลากหลายของมายกริดและทาเวอร์นา บริการที่นำเสนอโดยโครงการมายกริดแสดงด้วยบล็อกสี่เหลี่ยมเชื่อมส่วนบริการอื่นๆซึ่งอยู่กระจัดกระจายบนเครือข่ายอินเทอร์เน็ต เป็นบริการที่โครงการมายกริดมีข้อตกลงร่วมกันในการสื่อสารซึ่งกันและกันได้ บริการต่างๆเหล่านี้สามารถเข้าถึงได้โดยสาธารณะด้วยเครื่องมือหรือการโปรแกรมใดๆก็ได้ และเป็นบริการที่ได้รับการเชื่อมประสานเข้ากับโปรแกรมทาเวอร์นาแล้ว ซึ่งอาจมองว่าบริการภายนอกเหล่านี้เป็นปลั๊กอินหรือตัวเสริมโปรแกรมในโปรแกรมทาเวอร์นา

ในงานวิทยานิพนธ์นี้ได้สร้างบริการหรือเว็บเซอร์วิสขึ้นมาใช้งานเอง ในการทำนายโครงสร้างต้นไม้สายวิวัฒนาการของกิ้ง ซึ่งเป็นงานวิจัยร่วมกับศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศ คณะวิทยาศาสตร์ มหาวิทยาลัยสงขลานครินทร์ โดยทำงานอยู่บนระบบคลัสเตอร์ของศูนย์กริดแห่งมหาวิทยาลัยสงขลานครินทร์ ระบบคลัสเตอร์ติดตั้งระบบ Rocks Cluster รุ่น 4.3 [35] และสร้างส่วนประสานการเชื่อมต่อของเว็บเซอร์วิส ที่สร้างขึ้นมาใช้งานเองกับโปรแกรมทาเวอร์

นาโดยใช้โปรแกรม Soablab Analysis Tool [37] ซึ่งแสดงด้วยบล็อกสี่เหลี่ยมในไดอะแกรม นอกจากนี้ยังเลือกใช้บริการอื่นๆที่จำเป็น และมีการทำงานที่ถูกต้องตรงตามความต้องการในการ ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของกิ่งจากบริการของสถาบันชีวสารสนเทศแห่งยุโรป หรือ European Bioinformatics Institute (EBI) [39] ด้วย

เวิร์คโฟลว์ต่างๆในรูปแบบของเอกสาร Simple Conceptual Unified Flow Language (Scufl) [9] ที่สร้างในงานวิทยานิพนธ์แสดงด้วยบล็อกสี่เหลี่ยม ซึ่งเมื่อพิจารณาในเชิง สถาปัตยกรรมแล้ว เวิร์คโฟลว์ก็เป็นส่วนหนึ่งของการทำงานภายนอกหรือ Open world ก็ได้ ใน วิทยานิพนธ์นี้ได้เรียกใช้บริการชนิดหลักๆ คือ เว็บเซอร์วิสมาตรฐานแสดงด้วยบล็อกสี่เหลี่ยม ตัวอย่างเช่น Web API for Biology (WABI) ของ DNA Data Bank of Japan (DDBJ) [40] และ บริการไปโอมาร์ทแสดงด้วยบล็อกสี่เหลี่ยม โดยเวิร์คโฟลว์เรียกใช้งานบริการต่างๆผ่านแกนบริการ กลางของมายกริด (myGrid Core Services) ไปยังเว็บเซอร์วิสหรือบริการต่างๆที่กระจายกัน อยู่บนอินเทอร์เน็ต



รูปที่ 2.3 สถาปัตยกรรมของระบบที่เกี่ยวข้องในงานวิทยานิพนธ์ [9]

2.2 โปรแกรมทาวเวอร์นาและเวิร์คโฟลว์

ต่อไปจะกล่าวถึงภาพรวมและส่วนประกอบของโปรแกรมทาวเวอร์นาดังนี้

2.2.1 ภาพรวมของโปรแกรมทาวเวอร์นา

โปรแกรมทาวเวอร์นาอำนวยความสะดวกให้ผู้ใช้งาน สามารถสร้างเวิร์คโฟลว์ในการวิเคราะห์ระบบงานที่ซับซ้อนได้จากคอมโพเนนต์หรือบริการต่างๆ ทั้งที่อยู่ในท้องถิ่น หรือกระจัดกระจายอยู่บนอินเทอร์เน็ต แล้วสั่งให้ประมวลผลเพื่อให้ได้ผลลัพธ์ในรูปแบบที่เข้าใจได้ง่ายในการสนับสนุนกลไกหน้าที่การทำงานเหล่านี้ โปรแกรมทาวเวอร์นาจึงมีคอมโพเนนต์ต่างๆ มารองรับการทำงาน อย่างเช่นการค้นหา, การให้คำอธิบาย, และการเลือกบริการต่างๆ ที่เกี่ยวข้องกับการทำงานตามความต้องการนั้นๆ [2]

1) คำนิยาม

คำนิยามของคำศัพท์ที่ใช้ในสภาพแวดล้อมของทาวเวอร์นาในวิทยานิพนธ์มีดังนี้

- **เวิร์คโฟลว์ (Workflow)** – เป็นกลุ่มของคอมโพเนนต์และความสัมพันธ์ต่างๆ ที่กำหนดหรือบอกถึงกระบวนการที่ซับซ้อนกระบวนการหนึ่งจากการแสดงผลด้วยบล็อกอย่างง่าย ความสัมพันธ์ต่างๆนี้อาจจะอยู่ในรูปแบบของการเชื่อมโยงของข้อมูล (Data links), การส่งต่อข้อมูลซึ่งเป็นเอาต์พุตของคอมโพเนนต์หนึ่งๆ ไปเป็นอินพุตให้กับคอมโพเนนต์อื่นๆ (Data flow) หรืออยู่ในรูปแบบของการควบคุมการเชื่อมต่อ ที่บอกเงื่อนไขการทำงานของคอมโพเนนต์นั้นๆ (Control links) ตัวอย่างเช่น คอมโพเนนต์นี้จะทำงานก็ต่อเมื่อคอมโพเนนต์ที่ทำงานเสร็จแล้ว ในสิ่งแวดล้อมของทาวเวอร์นา เวิร์คโฟลว์หนึ่งๆ แสดงรูปแบบการไหลของข้อมูลหรือ Workflow data model ซึ่งอยู่ที่เครื่องส่วนบุคคลหรือบนเว็บไซต์ก็ได้โดยเป็นไฟล์ XML ในรูปแบบของเอกสาร Simple Conceptual Unified Flow Language (Scufl)

- **คอมโพเนนต์ (Component)** – คอมโพเนนต์จะใช้ในการสร้างบล็อกที่ทำหน้าที่อย่างใดอย่างหนึ่งซึ่งกำหนดไว้แล้วภายในกระบวนการนั้นๆ เราสามารถมองได้ว่าเครื่องมือรับคำสั่ง (Command line tool) ใดๆ หรือสคริปต์ใดๆ ของภาษา Perl ก็เป็นหนึ่งในคอมโพเนนต์ คอมโพเนนต์เหล่านี้ควรมีลักษณะการทำงานแบบอะตอมมิกหรือปรมาณู (Atomic) อธิบายลักษณะเช่น

ไม่สามารถแบ่งย่อยได้มากกว่านี้อีกแล้ว คอมโพเนนท์อาจจะเป็นหน่วยที่ทั้งใช้ข้อมูลและผลิตข้อมูลก็ได้ ตัวอย่างเช่น งานของโปรแกรม Basic Local Alignment Search Tool (BLAST) [41] หนึ่งๆ คือหนึ่งคอมโพเนนท์ที่ใช้ข้อมูลลำดับเบสและพารามิเตอร์สำหรับการค้นหาอื่นๆ และผลลัพธ์ก็คือผลิตรายงานที่บอกว่าลำดับเบสใดพบในฐานข้อมูลใดบ้าง และมีรายละเอียดอย่างไร เป็นต้น คอมโพเนนท์เหล่านี้อาจติดตั้งอยู่บนระบบที่มีพลังการคำนวณในระดับสูง หรือแม้แต่ในคอมพิวเตอร์ส่วนบุคคลของผู้ใช้งานเอง โดยสามารถเข้าถึงได้ง่ายผ่านเครือข่ายอินเทอร์เน็ต

- **บริการ (Service)** – บริการอีกชื่อหนึ่งก็คือคอมโพเนนท์ แต่ใช้คำว่าบริการในกรอบอ้างอิงว่า บริการเหล่านั้นทำงานอยู่บนเครื่องบริการหรือเซิร์ฟเวอร์ที่มีกำลังการประมวลผลสูงซึ่งไม่ได้เป็นเครื่องในท้องถิ่นของผู้ใช้เอง หรือบางกรณีก็กล่าวถึงบริการที่ปรากฏในเวิร์คโฟลว์ของโปรแกรมทาวเวอร์นา อาจจะใช้คำว่า โปรเซสเซอร์ (Processor) ก็ได้

- **เอ็นเนกเตอร์ (Enactor)** – เอ็นเนกเตอร์ของเวิร์คโฟลว์คือหน่วยที่ทำหน้าที่ในการประสานงานการทำงานของคอมโพเนนท์ต่างๆ ในเวิร์คโฟลว์ โดยอาจจะเรียกได้ว่าเป็นบริการๆหนึ่งโดยตัวมันเองในกรณีที่มีผู้ใช้เอกสารเวิร์คโฟลว์ และพารามิเตอร์ต่างๆ เป็น อินพุต และผลิตผลลัพธ์ออกมา หรือในกรณีของทาวเวอร์นาอาจเป็นคอมโพเนนท์หนึ่งภายในชุดของโปรแกรม เอ็นเนกเตอร์จะจัดการกระบวนการในภาพรวมคือการรายงานสถานะการทำงาน, การส่งผ่านข้อมูลระหว่างคอมโพเนนท์และงานอื่นๆที่เกี่ยวข้องทั้งหมด

2) ประโยชน์ของการใช้เทคโนโลยีเวิร์คโฟลว์

โปรแกรมทาวเวอร์นาสร้างขึ้นด้วยเทคโนโลยีใหม่ทั้งหมดเช่น Web service และ Ontology เป็นต้น โดยเกิดจากการพยายามลดอุปสรรคในการใช้งานของผู้ใช้ที่ต้องเสียเวลาไปมากกับการทำงานวิจัยบนฐานข้อมูลที่ซับซ้อนที่กระจายอยู่บนอินเทอร์เน็ต ประโยชน์ของการใช้เทคโนโลยีเวิร์คโฟลว์ มีดังนี้ [8][9]

- **เพิ่มประสิทธิภาพ (Efficiency)** – โปรแกรมทาวเวอร์นาช่วยประหยัดเวลาของผู้ใช้งานในหลายๆกรณีด้วยเทคโนโลยีทางด้านเวิร์คโฟลว์

- **ช่วยในการออกแบบการวิเคราะห์ (Analysis Design)** – การออกแบบการวิเคราะห์แบบเวิร์คโฟลว์จะสามารถทำได้รวดเร็วกว่าวิธีอื่นๆที่ผ่านมา ทำให้มองเห็นการวิเคราะห์ในภาพรวมได้ และยังสามารถเพิ่มคอมโพเนนท์ใหม่ๆเพื่อเป็นการขยายการวิเคราะห์ได้โดยง่าย

- **การเรียกใช้การทดลองแบบ *in silico* จากระยะไกล (Experiment Invocation)** – การทำงานงานด้านชีวสารสนเทศแบบดั้งเดิมส่วนมาก ตั้งแต่ขนาดเล็กรวมไปถึงโครงการการทำเครื่องหมายจีโนม (Whole genome annotation) ทำโดยการผสมผสานกันระหว่างเว็บเบราว์เซอร์และการใช้งานคำสั่งบนระบบปฏิบัติการ UNIX หรือเครื่องมือการทำงานอื่นๆ การผสมผสานของเครื่องมือเหล่านี้ ต้องการส่งผ่านข้อมูลระหว่างคอมพิวเตอร์ที่ต่างๆ โดยการคัดลอกแล้ววาง ลักษณะการทำงานแบบนี้ใช้เวลามาก มีแนวโน้มที่จะเกิดข้อผิดพลาดได้ง่าย และไม่เป็นระบบการทำงานที่ดี โปรแกรมทาวเวอร์น่ามีความสามารถในการรวบรวมคอมพิวเตอร์ที่ต่างๆเหล่านี้เข้าด้วยกันเป็นเวิร์คโพล์และเรียกใช้งานจากระยะไกลได้

- **จัดการคอมพิวเตอร์ (Component Management)** – ในยุคปัจจุบันมีความจำเป็นต้องเรียกใช้งานบริการฐานข้อมูลจากแหล่งระยะไกล เช่น สถาบันชีวสารสนเทศแห่งยุโรป หรือ European Bioinformatics Institute (EBI) ทำให้ผู้ใช้สามารถใช้งานบริการที่ทันสมัยอยู่เสมอๆ มีโปรแกรมหรือคอมพิวเตอร์ที่ต่างๆให้เลือกใช้ได้มากมาย และที่สำคัญคอมพิวเตอร์เหล่านั้นทำงานบนระบบคลัสเตอร์ โปรแกรมทาวเวอร์น่ามีความสามารถที่เข้าถึงเครื่องซูเปอร์คอมพิวเตอร์เหล่านี้ได้โดยง่ายผ่านทางเครื่องคอมพิวเตอร์ส่วนบุคคลสมัยใหม่ เนื่องจากโครงการมายกริดมีความร่วมมือกับแหล่งฐานข้อมูลมาตรฐานระดับโลกต่างๆ

การใช้งานเทคโนโลยีเวิร์คโพล์ ทำให้ไม่ต้องพัฒนาระบบงานที่ไม่จำเป็น เช่น อัลกอริทึมการทำนายโครงสร้างของโปรตีนขั้นที่สอง ก็จะต้องพัฒนาบริการที่เกี่ยวข้องทั้งหมดด้วยตัวเองเพื่อให้เพียงพอที่จะสามารถทำงานได้ ต้องสืบค้นข้อมูลลำดับเบสก่อน ซึ่งต้องมีการจัดการฐานข้อมูลเข้ามาเกี่ยวข้องโดยเป็นการเพิ่มภาระงานโดยไม่จำเป็น แต่หากใช้เทคโนโลยีเวิร์คโพล์ ผู้ใช้สามารถเลือกบริการที่พร้อมทำงานโดยผู้ใช้ไม่จำเป็นต้องสนใจเรื่องการเข้าถึงฐานข้อมูลเลย ผู้ใช้งานสามารถเลือกแหล่งบริการที่เหมาะสมตามความต้องการซึ่งมีให้เลือกมากมาย เช่นผู้ให้บริการหลักทางด้านชีวสารสนเทศอย่างสถาบันชีวสารสนเทศแห่งยุโรป (EBI) หรือไบโอมาร์ท (BioMart)

- **การเข้าถึงระบบประมวลผลที่มีประสิทธิภาพสูง (Invocation Performance)** – แม้ว่าเครื่องคอมพิวเตอร์บุคคลสมัยใหม่จะมีฮาร์ดแวร์ที่ดีขึ้นตามลำดับ แต่อย่างไรก็ตามก็ยังต้องการอัลกอริทึมในการทำงานเพื่อดึงเอาพลังการประมวลผลจากฮาร์ดแวร์เหล่านั้นมาใช้ การใช้งานคอมพิวเตอร์ในระยะเวลาไกล สามารถทำให้ผู้ใช้งานสามารถใช้ประโยชน์จากฮาร์ดแวร์ที่ทรงประสิทธิภาพ ตัวอย่างเช่นโปรแกรม InterProScan [42] ที่พัฒนาโดย EBI ไม่ได้ทำงานอยู่บนเครื่องคอมพิวเตอร์ธรรมดา แต่ทำงานบนเครื่องคลัสเตอร์ขนาดหลายร้อยโหนด ซึ่งเวิร์คโพล์ที่มีบริการ

ที่มีโครงสร้างพื้นฐานที่สนับสนุนการใช้บริการ จากแหล่งที่มีประสิทธิภาพการประมวลผลสูงแบบนี้ สามารถให้ผลลัพธ์แก่นักวิจัยที่รวดเร็วกว่าการใช้เครื่องคอมพิวเตอร์ธรรมดาอย่างมาก

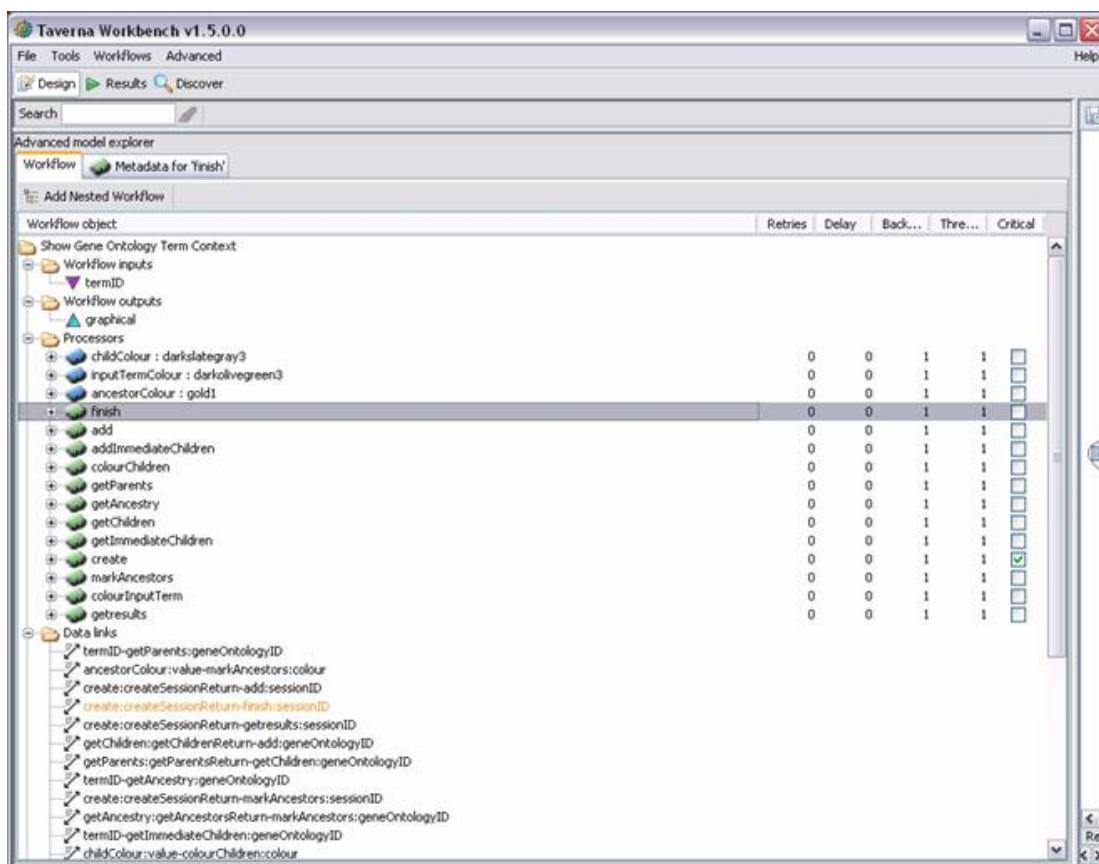
2.2.2 ส่วนประกอบของโปรแกรมทาวเวอร์นา

ในส่วนนี้จะนำเสนอส่วนประกอบทางด้านกราฟิกต่างๆ ของโปรแกรมทาวเวอร์นา และอธิบายการใช้งานเครื่องมือต่างๆ ในขั้นต้น เพื่อให้นำไปสู่ความเข้าใจในการออกแบบเวิร์คโฟลว์สำหรับแก้ปัญหาใดๆ ที่สนใจในวิทยานิพนธ์ ส่วนประกอบหลักๆ ของโปรแกรมทาวเวอร์นามีดังนี้

1) Advanced Model Explorer (AME)

- Entity Table

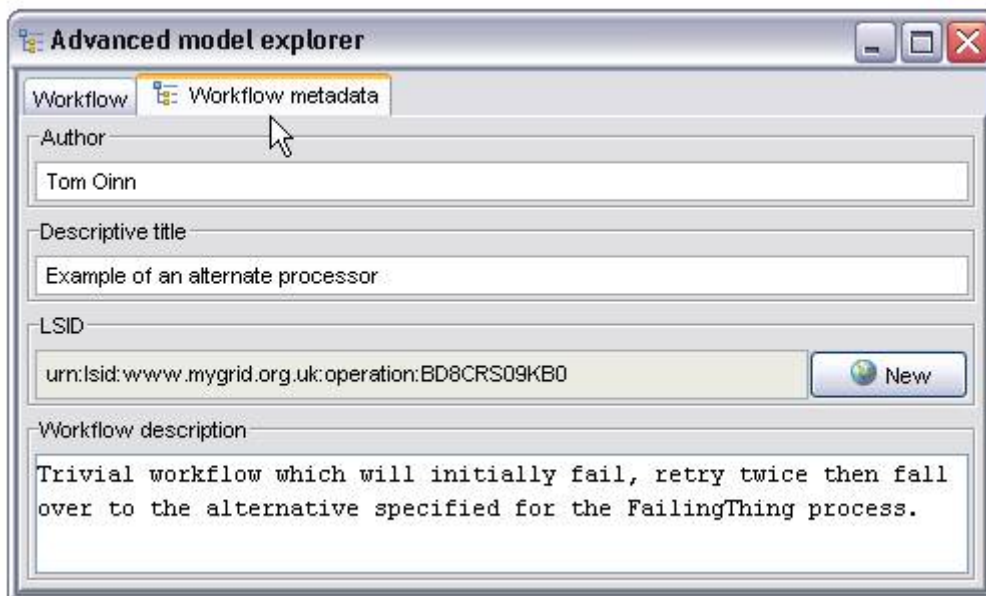
Entity Table เป็นส่วนหลักของ Advanced Model Explorer หรือ AME ซึ่งแสดงโปรเซสเซอร์หรือบริการ, อินพุตและเอาต์พุตของเวิร์คโฟลว์, การเชื่อมต่อของข้อมูล (Data flow link) และการประสานการเชื่อมต่อทั้งหมด โดยส่วนประกอบดังกล่าวนี้ อยู่ในรูปแบบ โครงสร้างต้นไม้ ผู้ใช้สามารถแตกโครงสร้างต้นไม้เพื่อดูรายละเอียดได้ ดังรูปที่ 2.4



รูปที่ 2.4 หน้าต่าง AME ของโปรแกรมทาวเวอร์นา

• Workflow Metadata

Workflow Metadata เป็นคำอธิบายการทำงานและจุดประสงค์ของเวิร์คโฟลว์ ทำให้มีการใช้เวิร์คโฟลว์ร่วมกันได้ คำอธิบายเวิร์คโฟลว์สามารถเพิ่มเติมหรือแก้ไขปรับปรุงได้ง่าย ตัวอย่างดังรูปที่ 2.5



รูปที่ 2.5 รายละเอียดของคำอธิบายใน AME

- **Workflow Inputs and Outputs**

เราสามารถสร้างพอร์ตของอินพุตและพอร์ตเอาต์พุต ได้ตามความต้องการที่จะรองรับข้อมูลจากบริการ ตัวอย่างเช่น บริการไปโอมาร์ทบางบริการให้อาท์พุตหรือแอดีตริบิวต์จำนวนมาก ทั้งนี้สามารถที่จะเลือกได้ว่า จะเอาข้อมูลเหล่านี้มาใช้ต่อหรือไม่

นอกจากนี้ยังสามารถกำหนดชนิดของข้อมูล (Data Type) ให้กับอินพุตหรือเอาท์พุตของเวิร์ค โฟลว์ได้ ซึ่งในโปรแกรมทาวเวอร์น่าจะเรียกว่า เมตาเดตา (Metadata) โดยสามารถที่จะกำหนดเพิ่มเติมได้ ซึ่งโปรแกรมทาวเวอร์น่าสามารถรองรับชนิดข้อมูลที่นำเสนอในตารางที่ 2.1

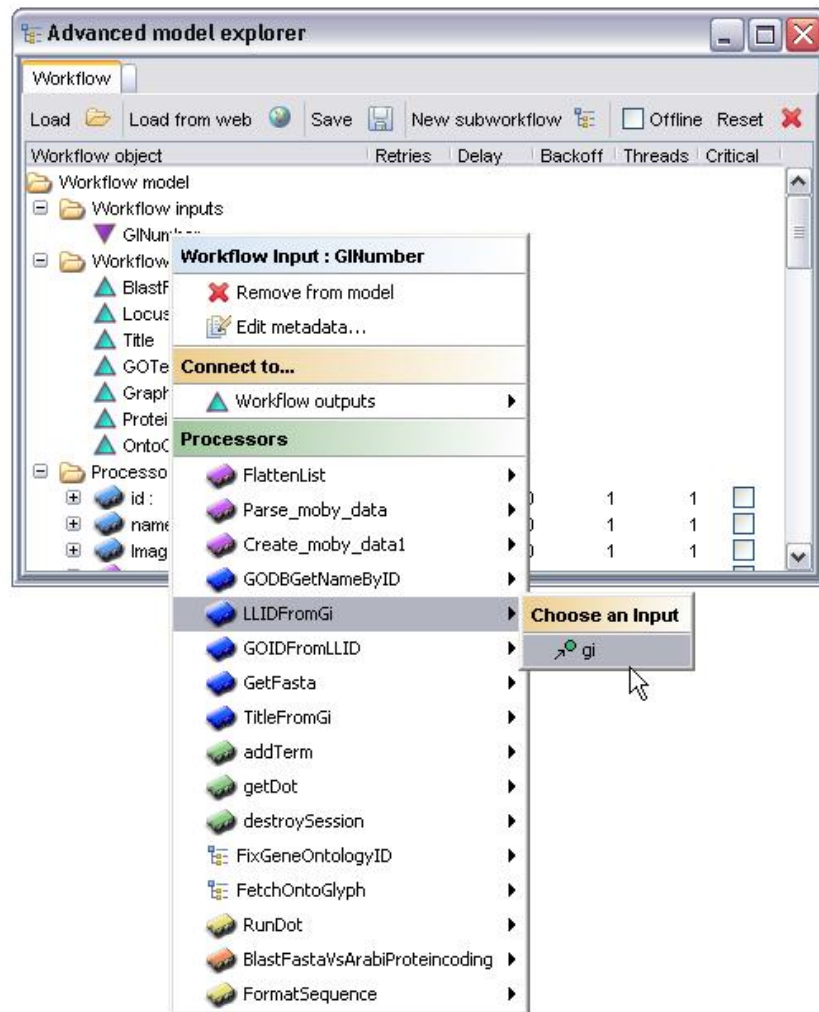
ตัวอย่างเช่น โปรแกรมทางด้านชีวสารสนเทศบางโปรแกรมให้อาท์พุตที่เป็นกราฟิกไฟล์ชนิด PNG ดังนั้น สามารถกำหนด MIME Type เป็น image/png=PNG Image ให้กับเอาท์พุตนั้นได้ เป็นต้น

ตารางที่ 2.1 ชนิดข้อมูลที่โปรแกรมทาวเวอร์นารองรับ

การกำหนดชนิดข้อมูล	ความหมาย
text/plain=Plain Text	ข้อความอักษรธรรมดา
text/xml=XML Text	ข้อความ XML
text/html=HTML Text	ข้อความ HTML
text/rtf=Rich Text Format	ข้อความอักษรธรรมดาพร้อมรูปแบบ
text/x-graphviz=Graphviz Dot File	ไฟล์ Dot ในระบบ Linux
image/png=PNG Image	ไฟล์ภาพสกุล PNG
image/jpeg=JPEG Image	ไฟล์ภาพสกุล JPEG
image/gif=GIF Image	ไฟล์ภาพสกุล GIF
application/zip=Zip File	ไฟล์บีบอัดสกุล ZIP
chemical/seq-aa-genpept=Genpept Protein	ข้อมูลลำดับของโปรตีน
chemical/seq-na-genbank=Genbank Nucleotide	ข้อมูลลำดับของนิวคลีโอไทด์
chemical/x-pdb=Protein Data Bank Flat File	ข้อมูลในรูปแบบโครงสร้างของโปรตีน

- การเชื่อมต่ออินพุตของเวิร์กโฟลว์เข้ากับโปรเซสเซอร์

โปรแกรมทาวเวอร์นารองรับจะแสดงพอร์ตของอินพุตของโปรเซสเซอร์ใดๆ ทุกพอร์ตสามารถระบุได้ว่า เอาต์พุตนี้เข้าไปเป็นพอร์ตอินพุตใดของโปรเซสเซอร์ใด เพราะโปรเซสเซอร์หรือบริการต่างๆ นั้นไม่ได้รับอินพุตหรือพารามิเตอร์เพียงตัวเดียว หากแต่สามารถรับอินพุตหรือพารามิเตอร์หลายๆ ตัวในการทำงานซึ่งขึ้นอยู่กับการใช้งานและการออกแบบ ดังรูปที่ 2.6



รูปที่ 2.6 การเชื่อมต่ออินพุตเข้ากับโปรเซสเซอร์

- **โปรเซสเซอร์ (Processor)**

โปรเซสเซอร์เป็นโหนดระดับบนสุดในโปรแกรมทาวเวอร์นา และเป็นคอมโพเนนต์ที่สำคัญที่สุด เพราะเป็นส่วนการทำงานของเวิร์กโฟลว์ ผู้ใช้งานสามารถรายละเอียดต่างๆ ของโปรเซสเซอร์ได้ ดังรูปที่ 2.7 ซึ่งแสดงผลสรุปให้ทราบว่าแต่ละโปรเซสเซอร์ว่าอยู่ที่ไหน (endpoint) เป็นบริการชนิดอะไร มีฟังก์ชันการทำงานอะไรบ้าง

Advanced model explorer

Workflow Remote resource usage

Save HTML description

Resource usage report

This display shows the various external resources used by the current workflow. It does not show resources such as local operations or string constants which are run within the enactment engine. Services are categorized by resource host and type, and the name of the instance of each service shown to the right.

Resources on mips.gsf.de , 1 instance.		
Biomoby	BlastFastaVsArabiProteincoding in /cgi-bin/proj/thal/contribute/MOBY/Service.cgi	BlastFastaVsArabiProteincoding
Resources on industry.ebi.ac.uk , 2 instances.		
Soaplab	Service rooted at /soap/soaplab	
	App category and name	Processors
	graphics::dot	RunDot
	edit::secret	FormatSequence
Resources on www.ebi.ac.uk , 4 instances.		
Web service	WSDL Defined at /collab/mygrid/service1/goviz/GoViz.jws?wsdl	
	Operation name	Processors
	createSession	createSession

รูปที่ 2.7 รายละเอียด Remote resource usage ของโปรเซสเซอร์ทั้งหมดที่ใช้ในเวิร์คโฟลว์

- การกำหนดค่าการคงทนต่อการทำงาน

การกำหนดค่าการคงทนต่อการทำงาน (Configuring Basic Fault Tolerance) ทำได้โดยการปรับแต่งค่าคุณสมบัติของแต่ละโปรเซสเซอร์ดังรูปที่ 2.8 อธิบายเฉพาะส่วนที่สำคัญดังนี้

- การทดลองใหม่ (Retries) การกำหนดจำนวนครั้งของทำงานซ้ำเมื่อการทำงานของโปรเซสเซอร์ล้มเหลว
- การถ่วงเวลา (Delay) ระหว่างการทดลองใหม่ในแต่ละครั้ง หน่วยเป็นมิลลิวินาที
- Backoff เป็นปัจจัยในการบอกว่า การทดลองใหม่ในแต่ละครั้งนั้น (ซึ่งจะอยู่ในช่วง n และ $n+1$) จะให้เพิ่มช่วงหน่วงเวลาเท่าไร คำนวณได้จาก delay time ยกกำลัง $n+1$ เช่น หากพิจารณาจากรูปที่ 2.6 กำหนดค่า delay time = 1 วินาที และ Backoff = 2 จะคำนวณได้ว่า

- รอบที่ 1 = $1000 \times 2^0 = 1000$ มิลลิวินาที = 1 วินาที
 - รอบที่ 2 = $1000 \times 2^1 = 2000$ มิลลิวินาที = 2 วินาที
 - รอบที่ 3 = $1000 \times 2^2 = 4000$ มิลลิวินาที = 4 วินาที
- ภาวะวิกฤต (Critical) เป็นการบอกว่า หากโปรเซสเซอร์ใดที่ถูกกำหนดให้มีภาวะวิกฤต หากโปรเซสเซอร์นั้นทำงานล้มเหลว ก็ให้หยุดการทำงานทั้งหมด แต่หากไม่ได้กำหนดไว้ หากโปรเซสเซอร์นั้นล้มเหลว เวอร์คโฟลว์ทั้งหมดจะยังคงทำงานได้ต่อไป แต่ส่วนที่เป็นดาวน์สตรีม (downstream) จากโปรเซสเซอร์ที่ล้มเหลวนี้จะไม่ถูกเรียกทำงาน



รูปที่ 2.8 กำหนดค่าการคงทนต่อการทำงาน

- โหนดแสดงการเชื่อมโยงของข้อมูล

โหนดแสดงการเชื่อมโยงของข้อมูล (Data Link Nodes) ที่ปรากฏใน AME คือ ข้อมูลแสดงการเชื่อมต่อของแต่ละโหนดในเวิร์คโฟลว์ โหนดเหล่านี้จะเกิดโดยอัตโนมัติเมื่อผู้ใช้ทำการต่ออินพุตและเอาต์พุตของโหนดที่เกี่ยวข้องในเวิร์คโฟลว์นั้นๆ การลบโหนดใดโหนดหนึ่งก็จะหมายถึงการยกเลิกการเชื่อมต่อของข้อมูลที่โหนดนั้นๆ

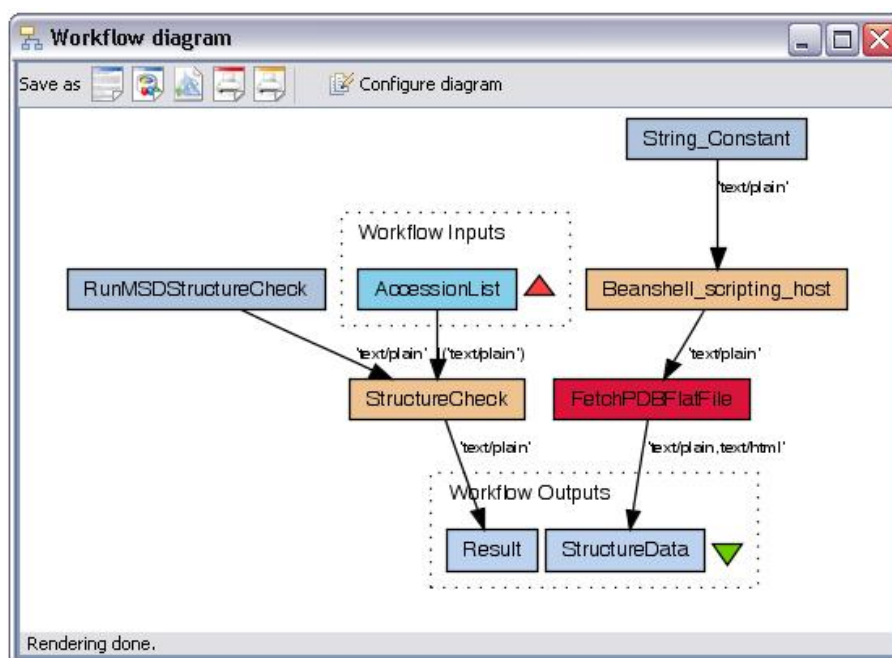
- โหนดแสดงการควบคุมการเชื่อมต่อ

โหนดแสดงการควบคุมการเชื่อมต่อ (Control link nodes) แสดงรายละเอียดการทำงานของโปรเซสเซอร์ที่ขึ้นอยู่กับเงื่อนไขนั้นๆ โหนดเหล่านี้จะเกิดโดยอัตโนมัติเมื่อผู้ใช้ทำการ

กำหนดเงื่อนไขต่างๆในเวิร์กโฟลว์ การลบโหนดใดโหนดหนึ่งก็จะหมายถึงการยกเลิกการควบคุมการเชื่อมต่อที่โหนดนั้นๆ

2) ไลออะแกรมของเวิร์กโฟลว์

ไลออะแกรมของเวิร์กโฟลว์ (Workflow diagram) เป็นส่วนประกอบที่ทำหน้าที่แสดงผลเพียงอย่างเดียวของเวิร์กโฟลว์ โดยอินพุต, เอาท์พุต, และโปรเซสเซอร์จะแสดงในรูปแบบของกล่องสี่เหลี่ยม และลูกศรแทนการเชื่อมต่อข้อมูล (Data links) และการควบคุมการเชื่อมต่อ (Control links) ดังรูปที่ 2.9



รูปที่ 2.9 ไลออะแกรมของเวิร์กโฟลว์

3) บริการที่พร้อมบริการ

เมื่อเริ่มต้นโปรแกรมทาวเวอร์นา โปรแกรมจะค้นหาบริการที่สามารถเข้าถึงได้และพร้อมให้บริการ (Available services) ทั้งหมดโดยโปรแกรมจะทำการตรวจสอบการตอบกลับ (Responding) ของเอ็นพอยต์ (Endpoint) ของบริการต่างๆด้วย หากบริการใดที่ไม่มีการตอบกลับก็จะไม่ปรากฏในรายชื่อบริการที่พร้อมให้บริการ ขณะเดียวกันเอ็นพอยต์ของบริการใดที่พร้อม

บริการหรือมีการตอบกลับ บริการเหล่านี้จะเป็นบริการที่ทันสมัยที่สุด ณ เวลาที่ผู้ใช้งาน โปรแกรม ตอนนั้นๆ ตัวอย่างดังรูปที่ 2.10 แสดง โหนดบริการต่างๆของ Soaplab จาก EBI และค่าเอ็นพอยต์ โดยปริยายกำหนดไว้ในไฟล์ *mygird.properties* โดยค่าโดยปริยายจะมีดังนี้

- เว็บเซอร์วิสที่อยู่บนพื้นฐานของมาตรฐาน Web Services Description Language (WSDL)
- Soaplab service คือ เว็บเซอร์วิสที่มีวิธีการการสร้างด้วยโปรแกรม Soaplab Analysis Tool Analysis Tool ซึ่งพัฒนาจาก EBI
- Biomoby คือเว็บเซอร์วิสจาก Biomoby
- BioMart คือเว็บเซอร์วิสจาก BioMart

นอกจากนี้ ผู้ใช้งานสามารถเพิ่มเอ็นพอยต์ที่ต้องการเข้ามาในโปรแกรมทาวเวอร์นา ได้ด้วยตนเองอีกด้วย เช่นในกรณีที่ผู้ใช้งานสร้างบริการขึ้นมาใช้งานเอง และสร้างระบบเชื่อมต่อ ประสานให้โปรแกรมทาวเวอร์นาเข้าถึงบริการนั้นได้ ซึ่งวิทยานิพนธ์นี้ก็จะนำเสนอวิธีการสร้างและ เข้าถึงบริการเหล่านี้เช่นกัน โดยจะกล่าวรายละเอียดต่อไป



รูปที่ 2.10 Available Service Panel แสดง โหนดของบริการ Soaplab จาก EBI

4) หน้าต่างแสดงผลการทำงาน

หน้าต่างแสดงผลการทำงาน (Result browser) ทำหน้าที่แสดงผลการทำงานเมื่อทุก โปรเซสเซอร์ในเวิร์คโพล์ทำงานเสร็จแล้ว หรือหากบางโปรเซสเซอร์ในเวิร์คโพล์ล้มเหลวใน

ระหว่างการทำงาน แต่ถ้ามีบางโปรเซสเซอร์ที่สามารถทำงานได้สำเร็จก็จะแสดงผลได้เช่นเดียวกัน นอกจากนี้ยังแสดงสถานภาพทำงานและรายงานโปรเซส (Process report) สำหรับคุณสถานะและรายงานสรุปการทำงานด้วย ดังรูปที่ 2.11 ผลลัพธ์เป็นกราฟิกชนิดข้อมูล PNG ซึ่งสามารถกำหนดชนิดของเอาต์พุตได้จาก AME ตาม MIME Type ที่โปรแกรมทาเวอร์น่าสนับสนุนไว้



รูปที่ 2.11 ผลลัพธ์การทำงานของเวิร์กโฟลว์จาก Result Browser

นอกจากนี้ทาเวอร์น่ายังมีตัวเรนเดอร์เอาต์พุต (Renderer) อื่นๆ ที่สนับสนุน อาทิ เช่น SeqVISTA [43] ซึ่งสนับสนุน MIME Type ดังนี้ Chemical/seq-aa-genpept, Chemical/seq-ngenbank, Chemical/x-swissprot, Chemical/x-embl-dl-nucleotide และ Chemical/x-ppd

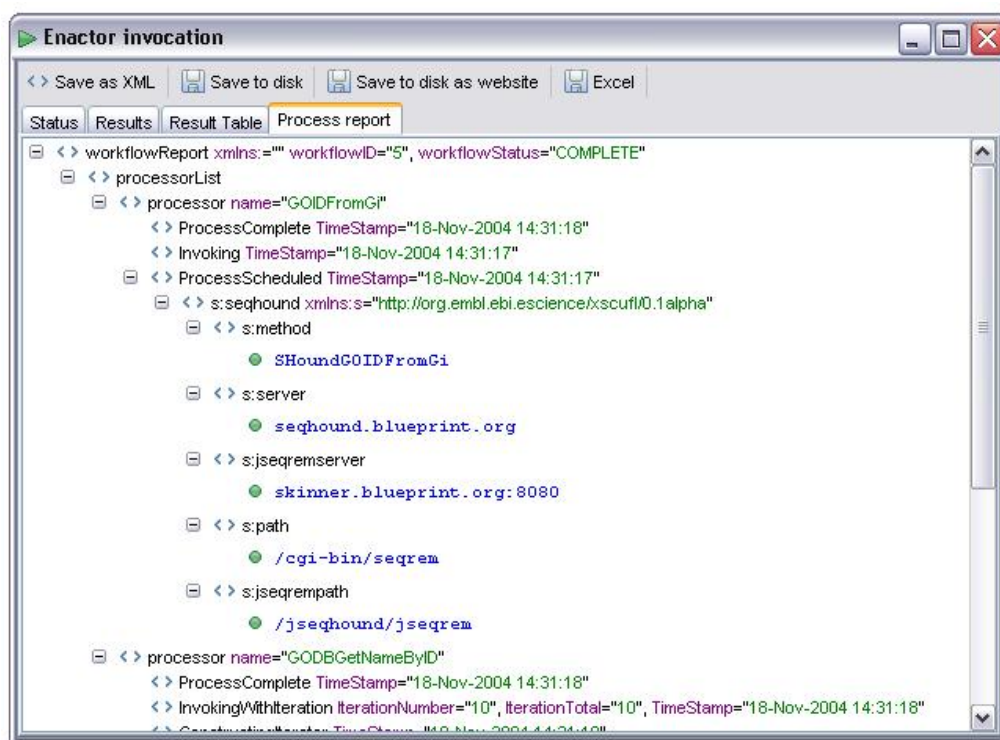
การบันทึกผลลัพธ์การทำงานในโปรแกรมทาเวอร์น่ามีหลายตัวเลือก ซึ่งแต่ละชนิดบันทึกผลลัพธ์จะแตกต่างกันไปตามลักษณะของรูปแบบที่เลือกดังนี้

- *Saving Entire Result Set as XML*: เป็นการบันทึกผลลัพธ์ในแต่ละรายการให้เป็นภาษา XML
- *Saving Entire Result Set as a Set of Files*: เป็นการบันทึกผลลัพธ์ทั้งหมดของ Workflow ลงดิสก์ ตามค่าชนิดของข้อมูลที่ได้กำหนดใน AME

- *Exporting Results to Excel*: เป็นการบันทึกข้อมูลชนิดข้อความ (Text) เป็นไฟล์ในรูปแบบของ Microsoft Excel ซึ่งจะช่วยในการจัดการข้อมูลต่อไป

- รายงานโปรเซส

รายงานโปรเซส (Process report) แสดงบันทึกการทำงานหรือล็อก (Log) การทำงานของเวิร์กโฟลว์เพื่อแสดงว่า เวิร์กโฟลว์ทำสำเร็จหรือไม่ อย่างไร โดยบันทึกรายละเอียดของทุกโปรเซสเซอร์ ซึ่งเป็นการช่วยในการดีบั๊ก (Debug) หรือตรวจสอบได้ทางหนึ่งว่าแต่ละโปรเซสเซอร์ว่าทำงานอย่างไร ใช้เวลาเท่าไร หรือล้มเหลวในช่วงไหนของการทำงาน ตัวอย่างดังรูปที่ 2.12



รูปที่ 2.12 รายละเอียดการทำงานของรายงานโปรเซส

2.3 Scuff และลักษณะเด่นของโปรแกรมทาวเวอร์นา

Scuff ย่อมาจาก Simple Conceptual Unified Flow Language [2][4] เป็นภาษาและมาตรฐานของโปรแกรมทาวเวอร์นาที่โครงการมายคริตสร้างขึ้นมา คุณสมบัติเด่นๆของภาษานี้มีดังนี้

2.3.1 การทำงานซ้ำๆ โดยปริยาย (Implicit iteration)

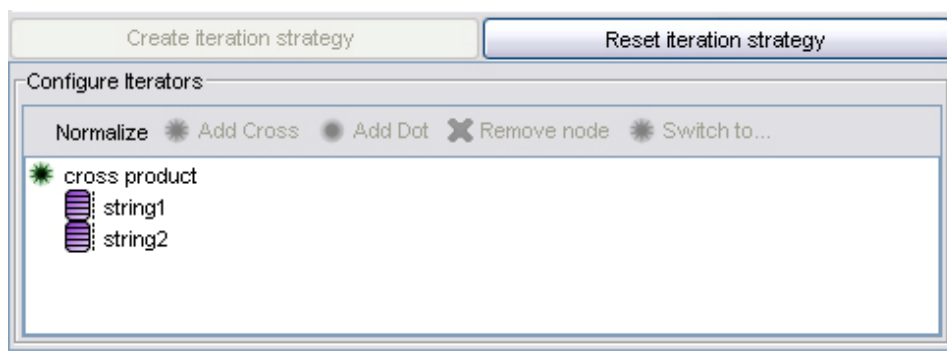
Iteration หรือการทำงานซ้ำ ในบริบทนี้คือการทำงานที่แบบซ้ำๆของกระบวนการใดการหนึ่งกับข้อมูลหลายรายการ ยกตัวอย่างเช่น กระบวนการทำงานการวิเคราะห์ Guanine-cytosine content (GC) จากเซตของลำดับพันธุกรรมหรือลำดับเบส ในภาษาการทำงานแบบดั้งเดิมอย่างเช่นภาษา Perl นั้น สามารถจัดการสิ่งเหล่านี้ได้ด้วยการทำเป็นลูป แล้วระบุเงื่อนไขการทำงานในลูปนั้นตัวอย่างเช่น `'while <condition> do <something>'` โดยการทำงานจะอยู่ภายใต้เงื่อนไขที่ระบุครบโดที่ยังคงมีรายการข้อมูลให้เหลือไว้ประมวลผล

สิ่งแวคล้อมของโปรแกรมทาวเวอร์นา นำเสนอกรอบการทำงานอย่างง่ายในการจัดการเกี่ยวกับลูป ผู้ใช้งานเพียงแต่ต่อเอาที่พุดซึ่งบรรจุเซตของรายการข้อมูลเข้ากับอินพุตซึ่งรับข้อมูลที่ละหนึ่งรายการได้เลย จากตัวอย่างการวิเคราะห์ GC ผู้ใช้อาจจะต้องการใช้บริการแก้หนึ่งบริการซึ่งใช้ข้อมูลลำดับเบส (DNA) หนึ่งลำดับ และผลิตเอาที่พุดเป็นเลขทศนิยมบอกถึงระดับความเข้มข้นของ GC ของลำดับนั้นๆ ในโปรแกรมทาวเวอร์นา ผู้ใช้สามารถต่อเซตของลำดับต่างๆ เข้าไปที่อินพุตเดียวกันได้เลย โดยโปรแกรมจะทำงานครั้งละหนึ่งลำดับจนข้อมูลอินพุตหมดสิ้น จากเอาที่พุดปกติคือเลขทศนิยมค่าเดียวต่อหนึ่งลำดับก็ไปเป็นรายการของเลขทศนิยม โดยรายการแรกของเอาที่พุดก็จะตรงกับรายการแรกของอินพุต เอาที่พุดที่สองก็จะตรงกับอินพุตที่สองแบบนี้ไปเรื่อยๆ

- **การทำงานซ้ำๆกับหลายอินพุต (Implicit iteration over multiple inputs)**

การทำงานโดยปริยายก็คือ ค่าต่างๆของอินพุตจะผสมผสานกันทั้งหมด จากรูปที่ 2.13 แสดงกลไกของการทำงานซ้ำๆกันของบริการหนึ่งแบบง่ายๆ โดยเป็นการเชื่อมต่อสายอักษร

สองสายเข้าด้วยกัน ในโปรแกรมทาเวอร์น่าสามารถเห็นการกำหนดค่าการซ้ำกันของข้อมูลโดยปริยายคือ หลายพอร์ตต่อหลายพอร์ต (All against all)

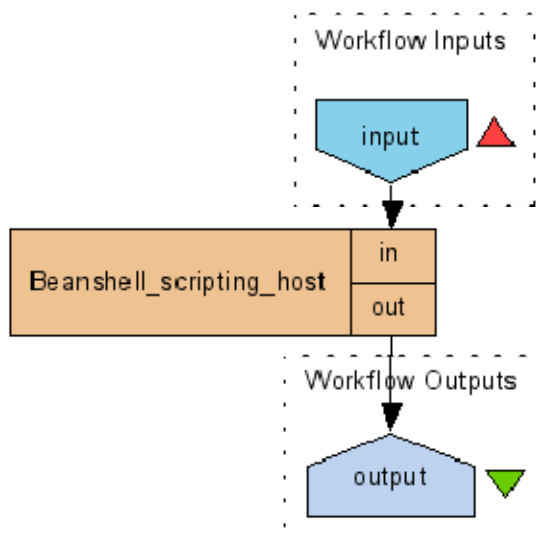


รูปที่ 2.13 Configure iterators ในโปรแกรมทาเวอร์น่า

กลไกของการซ้ำกันแสดงแบบโครงสร้างต้นไม้ โหนดใบของต้นไม้ตอบสนองต่ออินพุตนั้นๆ ให้กับบริการที่เชื่อมต่อเข้าด้วยกัน ในกรณีนี้คือ *String1* และ *String2* และที่ไม่ใช่โหนดข้อมูลคือวิธีการผสมผสานกันของอินพุต แบ่งเป็นสองชนิดคือ Cross product คือการผสมกันของพอร์ตข้อมูลแบบทั้งหมดต่อทั้งหมด (All against all) ซึ่งสามารถเปลี่ยนเป็น ลำดับแรกกับลำดับแรก ลำดับสองกับลำดับที่สอง แบบนี้ไปเรื่อยๆ และอย่างไรก็ตาม กลไกการปรับแต่งการซ้ำของข้อมูลสามารถกำหนด Cross product และ Dot product แบบผสมกันได้ตามลักษณะการทำงานของเวิร์คโฟลว์ ซึ่งในวิทยานิพนธ์จะใช้แบบผสมกัน

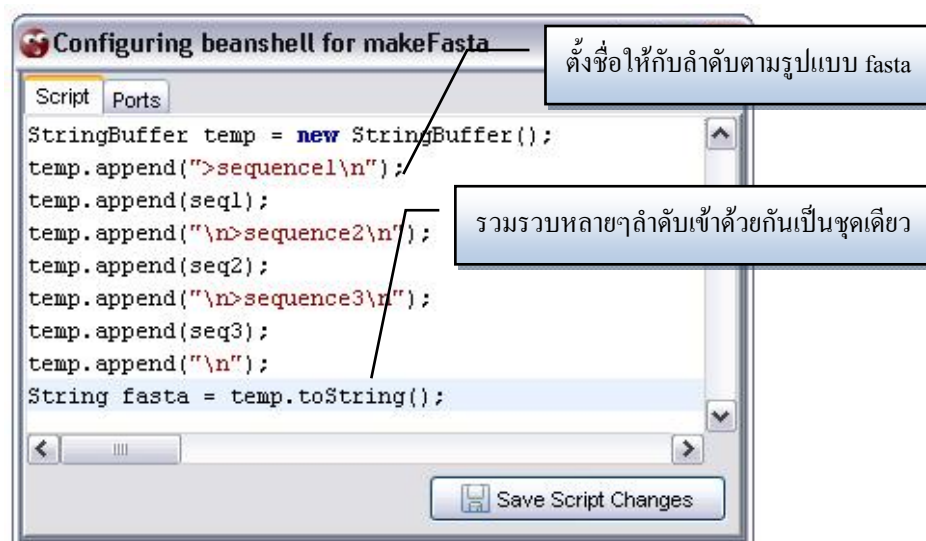
2.3.2 บีนเชลล์สคริปต์ (Beanshell scripting)

โดยปกติแล้วเอาต์พุตของบริการหนึ่งๆ อาจไม่ได้มีรูปแบบ (Format) ที่เหมาะสมและถูกต้องสำหรับบริการถัดไปหรือสำหรับ Nested workflow ได้ จึงต้องมีทางเลือกที่สามารถแก้ปัญหากรณีเหล่านี้ โดยบริการที่เข้ามาทำหน้าที่ดังกล่าวนี้เรียกว่า *Shim Service* [9] หรือเป็นบริการที่ไม่ซับซ้อนมาก ทำหน้าที่เปลี่ยนแปลงโครงสร้างหรือรูปแบบข้อมูลตลอดไปจนถึงการรองรับการเขียนภาษาโปรแกรมเพื่อรองรับอัลกอริทึมจากผู้ใช้งานเอง รูปที่ 2.14 แสดงเวิร์คโฟลว์การใช้งานบีนเชลล์อย่างง่ายโดยหมายถึงผลเอาต์พุตคืออินพุต (Output = Input)



รูปที่ 2.14 เวิร์กโฟลว์การใช้งานบีนเชลล์อย่างง่าย

บีนเชลล์สคริปต์ (Beanshell scripts) จัดเป็นบริการท้องถิ่นในโปรแกรมทาเวอร์น่า (Java local service) ที่ใช้ในการแก้ปัญหา โดยผู้ใช้สามารถเขียนภาษาสคริปต์เพื่อเข้าถึงตัวแปลคำสั่งของภาษา Java หรือ Java interpreter โดยไม่ต้องอาศัยความรู้พื้นฐานด้านการโปรแกรมภาษา Java มากนัก จากรูปที่ 2.15 เป็นตัวอย่างของสคริปต์ในการปรับแต่งสายอักขระของลำดับเบสซึ่งเป็นอินพุตให้กับบีนเชลล์ ซึ่งรวบรวมไปเป็นชุดข้อมูลเอาต์พุตแบบ Fasta (Fasta format) [36] โดยรวบรวมหลายๆลำดับเบสเข้าไว้ด้วยกัน ซึ่งในวิทยานิพนธ์นี้ใช้บีนเชลล์ในการสร้างสคริปต์เพื่อตรวจสอบบริการไบโอมาร์ท

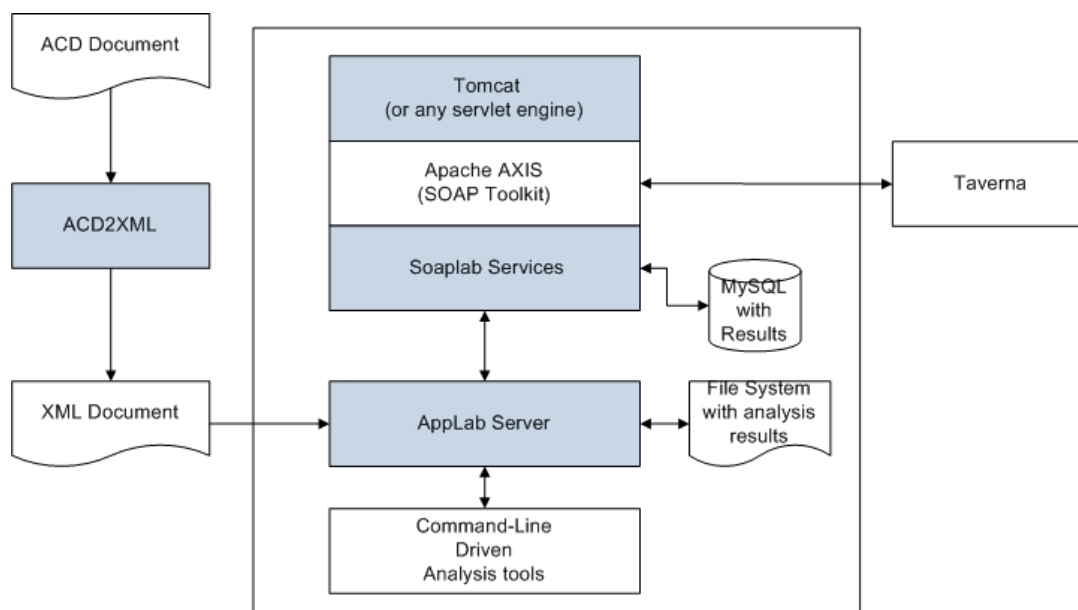


รูปที่ 2.15 ตัวอย่างสคริปต์ของบีนเชลล์

2.3.3 บริการ Soaplab

Soaplab Analysis Tool [37] คือกรอบการทำงานแบบหนึ่งที่สามารถทำให้โปรแกรมประยุกต์ที่อยู่ในรูปแบบของ Command-line สามารถทำงานเป็นเว็บเซอร์วิสได้ ส่วนที่สำคัญก็คือตัวระบบที่เป็นที่อยู่จัดการการทำงานของเว็บเซอร์วิส ในรูปที่ 2.16 คือสถาปัตยกรรมการทำงานของ Soaplab Analysis Tool จากรูป Tomcat คือระบบที่เป็นที่อยู่และตัวจัดการการทำงานของเว็บเซอร์วิส บริการ Soaplab ก็คือรายชื่อของเว็บเซอร์วิสต่างๆ ซึ่งมีหรือพอยน์เตอร์ (pointer) เชื่อมโยงไปยัง AppLab server ซึ่งเป็นระบบที่ครอบหรือห่อโปรแกรมประยุกต์ที่อยู่ในรูปแบบของ Command-line

Soaplab Analysis Tool เป็นคอมโพเนนท์ส่วนหนึ่งของมายกริดในโปรแกรมทาวเวอร์นา โดยมี SOAP Toolkit เป็นตัวเชื่อมประสานระหว่างกันซึ่งการติดต่อสื่อสารก็จะเป็นไปตามโปรโตคอล SOAP ในที่นี้คือ Apache AXIS [45] จะทำงานเป็น SOAP Toolkit



รูปที่ 2.16 สถาปัตยกรรมการทำงานของ Soaplab Analysis Tool [37]

2.3.4 การดีบั๊กและชี้นำทิศทางการทำงานของเวิร์คโฟลว์ด้วยเบรกพอยต์

ประโยชน์ของเวิร์คโฟลว์อีกอย่างหนึ่งคือ สามารถควบคุมทิศทางการทำงานได้ โดยในหัวข้อนี้จะอธิบายการชี้นำหรือการควบคุมทิศทางการทำงานของเวิร์คโฟลว์คือ การหยุดการทำงาน และแก้ค่าของตัวแปรต่างๆระหว่างทำงาน ซึ่งมีการติดต่อกับเวิร์คโฟลว์สองแบบคือ การหยุดชั่วคราว (Pausing) และการยกเลิก (Cancelling)

การหยุดการทำงานของเวิร์คโฟลว์สามารถทำได้สองแนวทางคือ โดยการใส่เบรกพอยต์ (Breakpoints) ที่จุดใดจุดหนึ่งของเวิร์คโฟลว์ซึ่งจะหยุดโดยอัตโนมัติ หรือจะหยุดการทำงานของเวิร์คโฟลว์เองทั้งหมดโดยการสั่งเองโดยผู้ใช้งาน (Manual) [2][4]

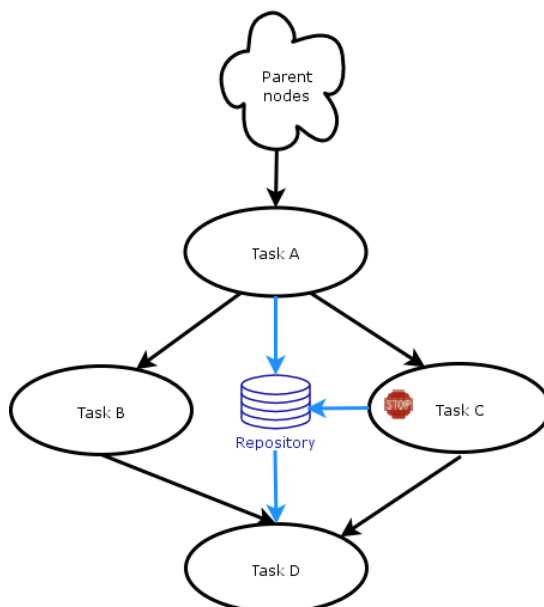
- **เบรกพอยต์ (Breakpoints)**

เบรกพอยต์เป็นศัพท์เทคนิคมาจากเครื่องมือการดีบั๊กซอฟต์แวร์ เครื่องมือเหล่านี้จะทำจุดหรือเครื่องหมาย (Mark) ในรหัสของโปรแกรม เพื่อให้โปรแกรมหยุดการทำงานให้ผู้ใช้งานสามารถเข้ามาแทรกแซงการทำงานของโปรแกรมได้ ในความหมายเดียวกัน เบรกพอยต์ใน

เวิร์คโฟลว์จะทำจุดหรือเครื่องหมายที่บริการหรือโปรเซสเซอร์ เพื่อให้หยุดการทำงาน ข้อมูลอินพุต และเอาต์พุตที่หยุดชั่วคราวที่บริการนั้นอาจจะแก้ไขหรือเปลี่ยนค่าได้

- **การแก้ไขข้อมูลระหว่างบริการหรือโปรเซสเซอร์**

การทำงานของเบรกพอยต์ในโปรแกรมทาวเวอร์นา จะตรงข้ามกับการดีบั๊กในวิศวกรรมซอฟต์แวร์ กล่าวคือเบรกพอยต์ในโปรแกรมทาวเวอร์นาจะหยุดการทำงานของบริการก็ต่อเมื่อบริการได้ผลิตเอาต์พุตเสร็จเรียบร้อยแล้ว ซึ่งเป็นการหยุดก่อนที่จะเริ่มต้นการทำงานของบริการถัดไป ซึ่งจุดนี้ผู้ใช้สามารถเข้าไปแก้ไขข้อมูลเหล่านี้ก่อนที่ส่งให้กับบริการถัดไปได้ เหตุผลของกระบวนการเหล่านี้คือ ต้องการที่จะรักษาความเป็นหนึ่งเดียวและความต่อเนื่องของข้อมูลในการกระจายและการทำงานแบบขนาน และยังป้องกันกรณีที่ไม่สามารถหยุดการทำงานของเวิร์คโฟลว์ได้ ยกตัวอย่างเช่น จากรูปที่ 2.17 บริการเอมีเอาต์พุตสองเอาต์พุตโดยแต่ละเอาต์พุตส่งต่อไปยังบริการบีและบริการซี ถ้าตัวดีบั๊กแบบต่างๆ ไปเพิ่มเบรกพอยต์ที่บริการซีและเปลี่ยนข้อมูลที่บริการเอ จะทำให้บริการซีมีข้อมูลใหม่ ในขณะที่บริการบียังมีข้อมูลเดิมอยู่ สถานการณ์แบบนี้เป็นการเพิ่มเงื่อนไข หรือมีแนวโน้มที่จะทำให้เวิร์คโฟลว์ผลิตผลลัพธ์ที่ไม่ถูกต้องหรือตรงตามความต้องการ นอกจากนี้ในขณะที่เวิร์คโฟลว์กำลังทำงาน อินพุตหรือเอาต์พุตจะถูกบันทึกไว้ใน Repository ซึ่งข้อมูลเดียวกันก็จะเก็บทั้งข้อมูลรุ่นเก่าและรุ่นใหม่ ดังนั้นในการแก้ไขหรือตรวจสอบข้อมูลเหล่านี้ที่จุดใดจุดหนึ่งของเวิร์คโฟลว์ ควรจะกำหนดเบรกพอยต์ที่บริการที่กำลังจะผลิตข้อมูล (หรือกำลังจะประมวลผลต่อ) และไม่ใช่บริการที่จะรับข้อมูลนั้นมาเป็นอินพุต



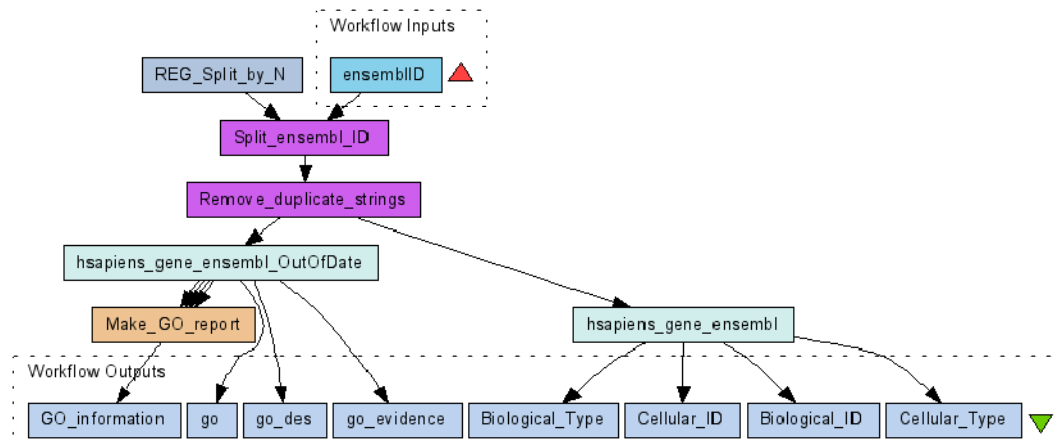
รูปที่ 2.17 ตัวอย่างกลไกการทำงานของเวิร์กโฟลว์เมื่อมีการกำหนดเอ็นพอยต์ [9]

2.4 ทาเวอร์น่าสอง รุ่นทดลองใช้งาน (Taverna 2 Preview)

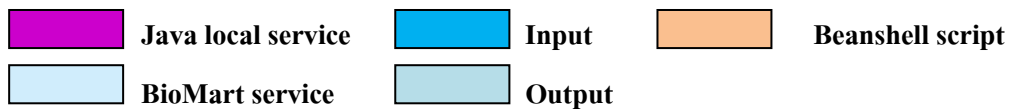
โปรแกรมทาเวอร์น่าตั้งแต่รุ่น 1.7.0 ก็มีปลั๊กอินชื่อว่า *Taverna 2 preview* ซึ่งเป็นการทดลองโปรแกรมทาเวอร์น่ารุ่น 2 เพื่อการทำงานที่เรียกว่า *Workflow health check* หรือการตรวจสอบเสถียรภาพของเวิร์กโฟลว์ [8][9] โดยจะแปลงตัวแบบเอกสาร Scufi แบบเก่าให้เป็นแบบ T2 style แล้วตรวจสอบว่าเอ็นพอยต์ของเว็บเซอร์วิสต่างๆ ในเวิร์กโฟลว์ตอบกลับผ่าน HTTP (Responding) หรือไม่ หากมีการตอบกลับก็คือสามารถเข้าถึงได้ หากไม่มีการตอบกลับ ปัญหาที่อาจจะมาจากเครือข่ายของผู้ใช้งานเองไม่สามารถติดต่อภายนอกได้ หรือไม่มีปัญหาที่ฝั่งของผู้ให้บริการเองอันเนื่องมาจากหลายๆสาเหตุ

ปลั๊กอิน *Taverna 2 preview* ตรวจสอบเฉพาะการตอบกลับ (Responding) ของเอ็นพอยต์ของบริการเท่านั้น ไม่ได้ตรวจสอบถึงการเปลี่ยนแปลงอันเนื่องมาจากการปรับปรุง (upgrade) ของบริการนั้นๆด้วย ซึ่งในงานวิจัยนี้เน้นความสนใจไปที่บริการไบโอมาร์ท ซึ่งมีการใช้งานทั้งฟิลเตอร์และแอ็คตริบิวต์ จากรูปที่ 2.18 แสดงเวิร์กโฟลว์ในการค้นหาข้อมูล Gene Ontology หรือโครงสร้างลำดับการให้คำอธิบายกระบวนการทางชีววิทยา [46] การระบุรายละเอียดในระดับเซลล์ (Cellular localization) และรวมถึงการทำงานในระดับโมเลกุล (Molecular function) โดยมีอินพุตคือ Ensembl IDs เวิร์กโฟลว์เรียกใช้บริการไบโอมาร์ทชื่อว่า *hsapiens_gene_ensembl* จำนวน

2 บริการ โดยบริการที่ได้สร้างก่อนหน้านี้และล้าสมัยไปแล้ว และได้ตั้งชื่อบริการที่ล้าสมัยว่า *hsapiens_gene_ensembl_OutOfDate* ซึ่งจะผลิตผลลัพธ์คือ Gene ontology, Gene ontology description และ Gene ontology evidence จากนั้นใช้สคริปต์ของบีนเชลล์ในการเก็บผลลัพธ์เหล่านั้น มารวบรวมเข้าไว้ด้วยกันเพื่อทำเป็นรายงาน

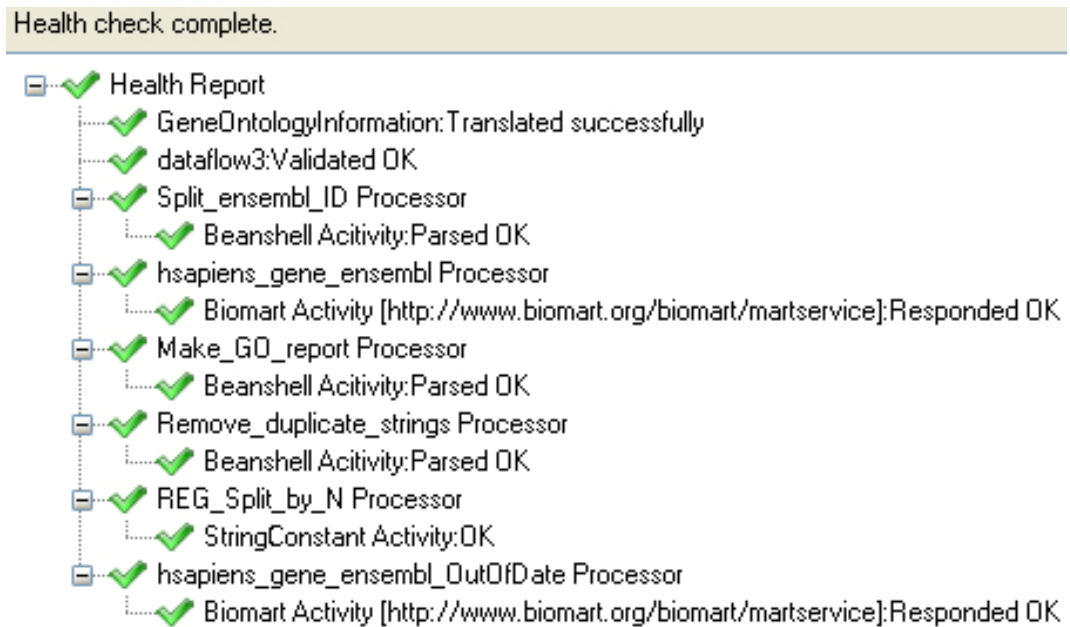


Legends:

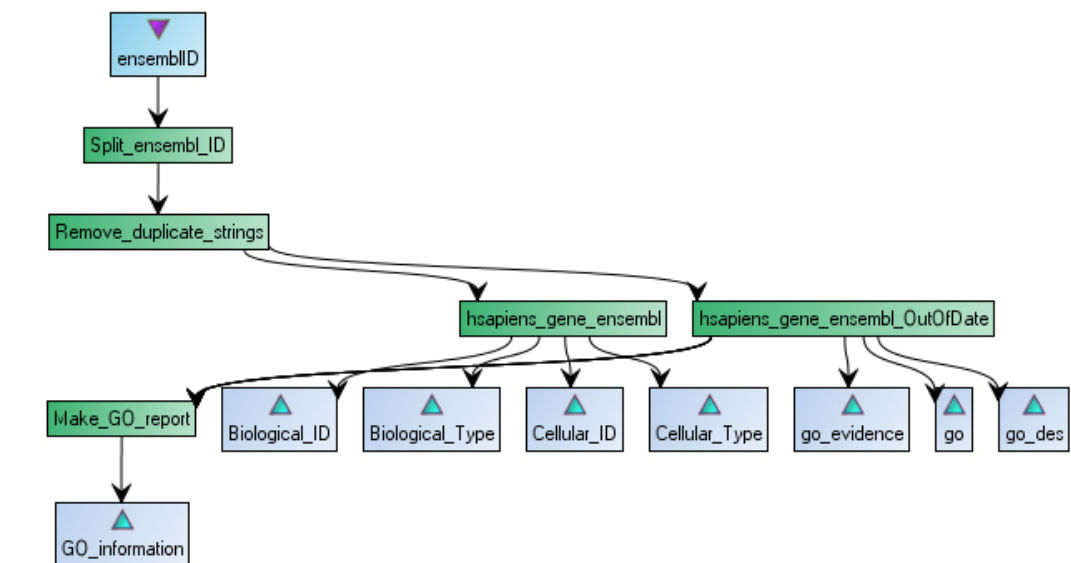


รูปที่ 2.18 เวิร์คโฟลว์ในการค้นหา Gene Ontology จาก Ensembl IDs

เวิร์คโฟลว์ในรูปที่ 2.18 ตรวจสอบด้วย *Taverna 2 preview* และได้ *Health check report* ดังรูปที่ 2.19 ซึ่งแสดงให้เห็นว่า ทุกบริการมีการตอบกลับทั้งสิ้น และในขั้นตอนต่อไปคือ เรียกการทำงานของเวิร์คโฟลว์ และได้สถานะการทำงานของเวิร์คโฟลว์ดังรูปที่ 2.20



รูปที่ 2.19 Health check report ของเวิร์คโฟลว์



รูปที่ 2.20 สถานะการทำงานของเวิร์คโฟลว์

จากรูปที่ 2.20 บริการสามารถทำงานได้สำเร็จโดยไม่มีข้อผิดพลาดใดๆ แต่จะพบว่า *hsapiens_gene_ensembl_OutOfDate* ไม่ได้ผลผลิตผลลัพธ์ใดๆเลย (empty result list) เนื่องจากแอ็ตทริบิวต์ที่ต้องการคือ Gene ontology, Gene ontology description และ Gene ontology

evidence ได้ล้มสมัยเสียแล้ว อย่างไรก็ตามเป็นเซลล์ก็ยังคงทำงานได้สำเร็จ เนื่องจากไม่ได้สนใจความมีอยู่ของข้อมูล (Data availability) เพียงแต่ตรวจสอบเฉพาะความผิดพลาดของไวยากรณ์ขณะที่ทำงานเท่านั้น ดังนั้นในวิทยานิพนธ์นี้จึงได้นำเสนอแนวทางแก้ไขปัญหาดังกล่าวโดยการพัฒนาเวิร์คโฟลว์ที่สามารถตรวจสอบฟิลเตอร์และแอ็คตริบิวต์ที่ล้มสมัยได้

2.5 บริการไบโอมาร์ท (BioMart)

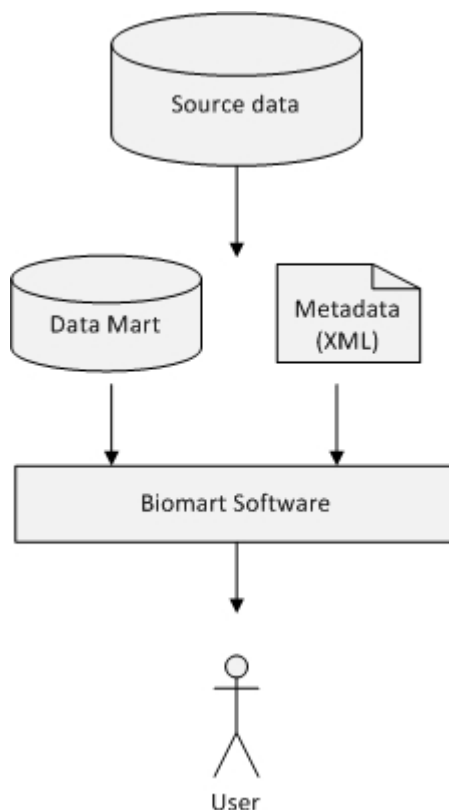
ต่อไปจะกล่าวแนะนำบริการไบโอมาร์ท คำอธิบายบริการ สถาปัตยกรรมการทำงานในภาพรวมของระบบบริการ และการเข้าถึงข้อมูลของบริการไบโอมาร์ทที่นำมาประยุกต์ใช้สำหรับงานวิทยานิพนธ์ มีดังนี้

2.5.1 แนะนำบริการไบโอมาร์ท

เว็บไซต์ของโครงการไบโอมาร์ทคือ <http://www.biomart.org> ซึ่งสามารถดาวน์โหลดโปรแกรมไบโอมาร์ทรุ่นล่าสุดและเอกสารคู่มือที่เกี่ยวข้องต่างๆได้ โครงการไบโอมาร์ทก่อตั้งร่วมกันโดยสถาบันชีวสารสนเทศแห่งยุโรป (EBI) และศูนย์วิจัยมะเร็ง Cold Spring Harbor [12]

2.5.2 คำอธิบาย

ไบโอมาร์ท คือระบบจัดการข้อมูลแบบ Open source โดยมีระบบเชื่อมประสานให้ผู้ใช้งานสามารถคิวรีข้อมูลชีวสารสนเทศได้หลากหลายและรวดเร็ว [11][12] จุดประสงค์ของไบโอมาร์ทคือ ให้บริการการแปลงข้อมูลใดๆไม่ว่าจะเป็นเอกสารข้อมูลต่างๆหรือฐานข้อมูลไปเป็นข้อมูลที่เรียกว่า ดาตามาร์ท (Data marts) ซึ่งสามารถเข้าถึงด้วยเว็บเบราว์เซอร์มาตรฐานและยังสามารถเชื่อมต่อกับภาษาโปรแกรม Perl, Java และ API ของเว็บเซอร์วิสอีกด้วย ดังรูปที่ 2.21



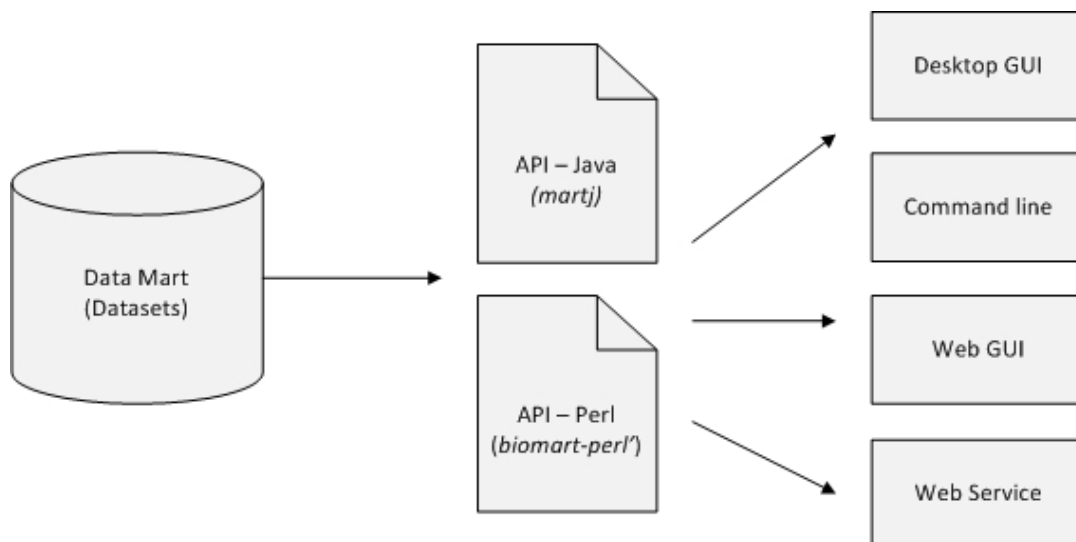
รูปที่ 2.21 ภาพรวมของไบโอมาร์ท [12]

ไบโอมาร์ทสร้างขึ้นเพื่อสนับสนุนการคิวรีข้อมูลได้อย่างดี ผู้ใช้งานสามารถเรียกใช้งานได้อย่างรวดเร็ว, ง่าย, และมีความสามารถในการคิวรีข้อมูลได้อย่างมีประสิทธิภาพ โดยการเลือกเงื่อนไขได้จากเว็บ, กราฟิก, และข้อความ การคิวรีข้อมูลทำได้โดยผ่าน API ของเว็บ เซอร์วิสหรือการติดต่อกับไลบรารีโดยใช้ภาษาโปรแกรม Perl และ Java และชุดโปรแกรมทั้งหมดของไบโอมาร์ทสามารถติดตั้งในเครื่องคอมพิวเตอร์ผู้ใช้งานเองได้ด้วย ไบโอมาร์ทเป็นระบบซอฟต์แวร์แบบเปิดภายใต้ลิขสิทธิ์ของ GNU Lesser General Public License (LGPL) [47] ซึ่งทุกคนสามารถเข้าถึงและใช้งานโดยไม่มีข้อผูกมัดหรือเงื่อนไขใดๆ

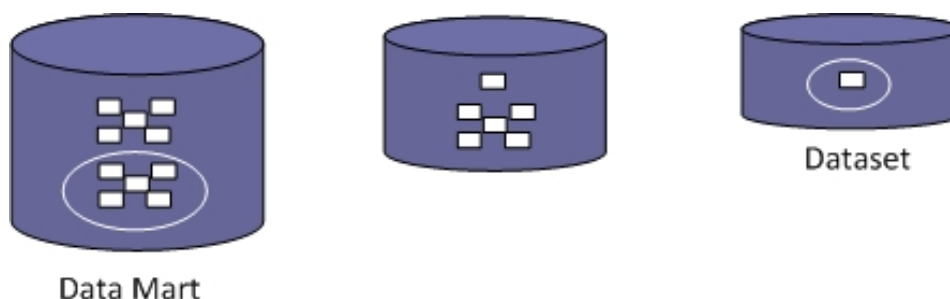
2.5.3 สถาปัตยกรรมของไบโอมาร์ท

ไบโอมาร์ทออกแบบโดยใช้สถาปัตยกรรมแบบ Three Tier Architecture ดังรูปที่ 2.22 โดยส่วนแรกประกอบด้วยฐานข้อมูลเชิงความสัมพันธ์หนึ่งฐานหรือมากกว่า ฐานข้อมูลแต่ละ

ฐานข้อมูลสามารถรองรับได้มากกว่าหนึ่งดาตามาร์ท และในแต่ละดาตามาร์ทจะมีดาตาเซต (Dataset) ได้หลายๆดาตาเซตดังรูปที่ 2.23 การกำหนดการปรับแต่งของดาตาเซตเหล่านี้ จะเก็บอยู่ตารางในแต่ละดาตามาร์ทนั้นๆ ซึ่งสามารถสร้างได้โดยใช้เครื่องมือที่เรียกว่า MartEditor



รูปที่ 2.22 สถาปัตยกรรมแบบ Three Tier Architecture ของไบโอมาร์ท [12]



รูปที่ 2.23 ความสัมพันธ์ระหว่างดาตามาร์ทและดาตาเซต [12]

ส่วนที่สองประกอบด้วย API จำนวนสองชนิดคือ ชนิดที่เขียนขึ้นด้วยภาษา Perl โดยกระจายอยู่ในไลบรารีชื่อ *biomart-perl* และอีกชนิดคือเขียนด้วยภาษา Java โดยกระจายอยู่ในไลบรารีชื่อ *martj*

ส่วนที่สามประกอบด้วยตัวเชื่อมประสานสำหรับงานคิวรีดังนี้

- MartView เป็นตัวเชื่อมประสานที่ทำงานกับเว็บเบราว์เซอร์ ทำงานบนพื้นฐานของ Perl

- MartService เป็นตัวเชื่อมประสานแบบเว็บเซอร์วิส ทำงานบนพื้นฐานของ Perl
- MartURLAccess เป็นตัวเชื่อมประสานของกลไกการเข้าถึง MartView ผ่านทาง URL ทำงานบนพื้นฐานของ Perl
- MartExplorer เป็นตัวเชื่อมประสานกับผู้ใช้แบบกราฟิก ทำงานบนพื้นฐานของ Java
- MartShell เป็นตัวเชื่อมประสานแบบเชลล์รับคำสั่ง ทำงานบนพื้นฐานของ Java

ข้อกำหนดการปรับแต่งของคาตาเซ็ทจะเก็บในรูปแบบของ XML ในตารางพิเศษของฐานข้อมูลที่คาตามาร์ทนั้นๆ ผู้ใช้จะจัดการเอกสารรีจิสทรี XML บนโคลเอ็นท์หนึ่งๆโดยจะระบุลงไปว่าคาตาเซ็ทอะไร ในคาตามาร์ทไหนและอยู่ในฐานข้อมูลอะไรที่พร้อมให้บริการสำหรับคิวรีข้อมูล ไปโอมาร์ทรุ่น 0.7 สนับสนุนระบบฐานข้อมูลเชิงสัมพันธ์อย่างเช่น MySQL, Oracle และ Postgres

2.5.4 การเข้าถึงข้อมูลของบริการไปโอมาร์ท

ในการพัฒนาเวิร์คโฟลว์สำหรับตรวจสอบบริการไปโอมาร์ทนี้ ได้ศึกษาถึงกลไกการเข้าถึงข้อมูลที่เกี่ยวข้องของบริการไปโอมาร์ทดังนี้

มาร์ทวิว (MartView) คือตัวเชื่อมประสานทำงานบนเว็บเบราว์เซอร์และคาตาเซ็ทของไปโอมาร์ทเพื่อให้ผู้ใช้งานสามารถติดต่อสื่อสารกันได้ โดยมาร์ทวิวถูกออกแบบให้ง่ายต่อการใช้งานที่สุดเท่าที่เป็นไปได้ มาร์ทวิวแบ่งตามการใช้งานได้ 2 ประเภทต่อไปนี้

1) การติดต่อด้วยเว็บเบราว์เซอร์ (Web browser interface)

มาร์ทวิวสามารถเข้าถึงได้ด้วยเว็บเบราว์เซอร์โดยมี URL ที่ขึ้นอยู่กับารปรับแต่งของผู้ดูแลระบบที่นำชุดโปรแกรมไปโอมาร์ทมาใช้งานในองค์กรนั้นๆโดยมี URL ศูนย์กลางของมาร์ทวิวที่ใช้ฐานข้อมูลกลางคือ <http://www.biomart.org/biomart/martview>

การใช้บริการไบโอมาร์ทด้วยเว็บเบราว์เซอร์ จะมีประโยชน์ในการทดสอบการทำงานเบื้องต้นได้ทันที และสามารถตรวจสอบได้ว่าให้ผลลัพธ์ที่ต้องการหรือไม่ นอกจากนี้ยังใช้วิธีการนี้ในการดีบั๊กหรือตรวจสอบความถูกต้องของบริการอีกด้วย เพราะบริการไบโอมาร์ทที่เข้าถึงผ่านเว็บเบราว์เซอร์ในขณะนั้นๆย่อมเป็นบริการที่ทันสมัยที่สุด

2) การติดต่อด้วยเว็บเซอร์วิส (Web Services API)

เอพีไอเว็บเซอร์วิสของไบโอมาร์ทมีชื่อว่า มาร์ทเซอร์วิส (MartService) โดยให้บริการเป็นส่วนหนึ่งมาร์ทวิว (MartView) สามารถเข้าถึงได้โดยการร้องขอผ่าน HTTP ซึ่งสามารถใช้เว็บเบราว์เซอร์มาตรฐานเข้าถึงได้ การเข้าถึงมาร์ทวิวโดยใช้ URL นี้

<http://www.mycompany.com/scripts/biomart/martview>

เข้าถึงเว็บเซอร์วิสได้ใช้ URL นี้

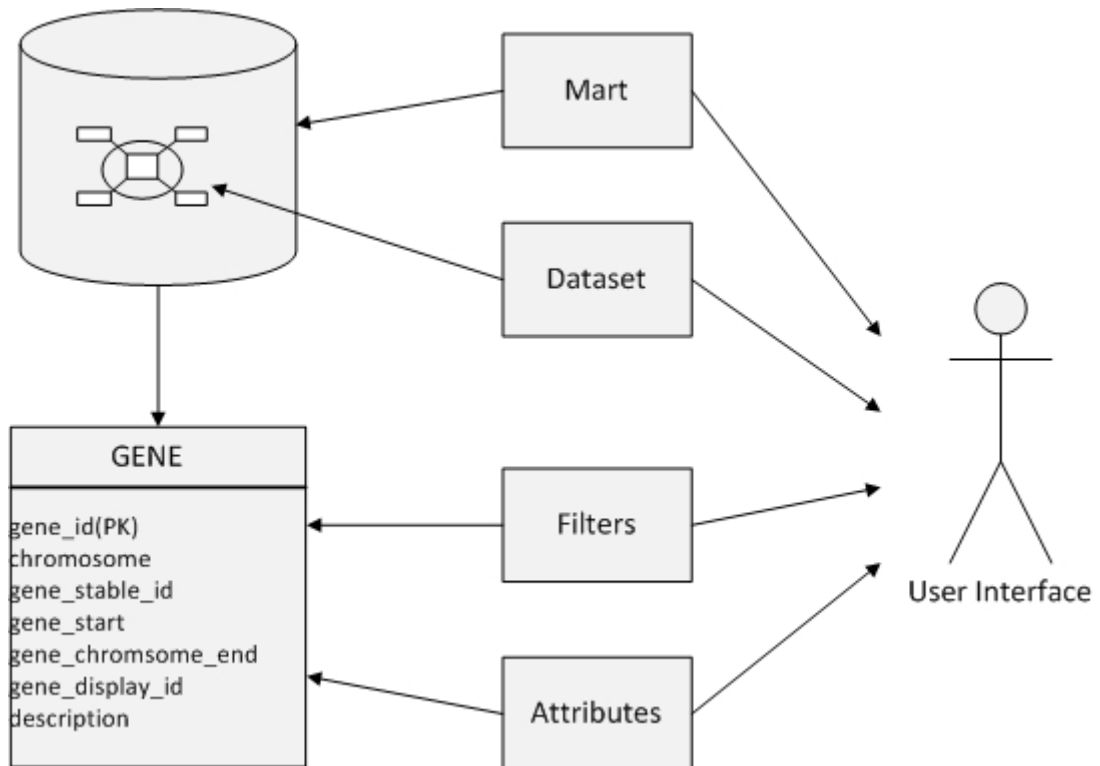
<http://www.mycompany.com/scripts/biomart/martservice>

และศูนย์กลางเว็บเซอร์วิสไบโอมาร์ทคือ

<http://www.biomart.org/biomart/martservice>

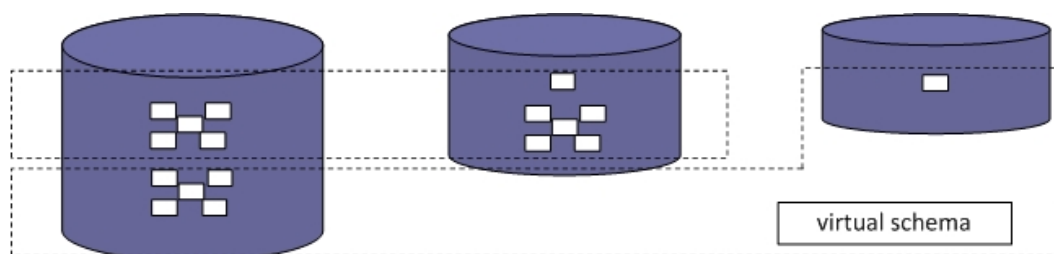
2.1) เมตาเดตา (Metadata)

เมตาเดตา (Metadata) หรือคำอธิบายข้อมูลของเว็บเซอร์วิสต่างๆ สามารถเข้าถึงได้ด้วยการร้องขอ 'GET' ในขณะที่คำสั่งคิวรีข้อมูลสามารถส่งได้ด้วยการร้องขอ 'POST' ของ HTTP โดยในดาตาเซตประกอบด้วยฟิลเตอร์หรือฟิลด์ข้อมูลที่ใช้เป็นเงื่อนไขการคิวรีและแอตทริบิวต์หรือฟิลด์ข้อมูลที่เป็นเอาต์พุตจากคิวรีนั้นๆ ตัวอย่างเช่นรูปที่ 2.24 ดาตาเซตของยีนซึ่งประกอบด้วยฟิลเตอร์และแอตทริบิวต์จำนวนมาก ยกตัวอย่างการคิวรีข้อมูลต่อไปนี้คือ หมายเลขยีน (Gene_stable_id), ตำแหน่งเริ่มต้น (Gene_start), ตำแหน่งจบ (Gene_chromosome_end), ชื่อของยีน (Gene_display_id) และคำอธิบาย (Description) โดยมีจะฟิลเตอร์เกี่ยวข้องคือ ระบุการค้นหาจากหมายเลขยีนที่ผู้ใช้เป็นผู้กำหนดเอง และระบุว่าจะค้นที่โครโมโซมไหน เช่น โครโมโซม 1 เป็นต้น



รูปที่ 2.24 ตัวอย่างเมตาดาตาที่ใช้ในการคิวรีข้อมูลจากไบโอมาร์ท [12]

ชุดเมตาดาตาของดาตาเซ็ทที่ให้บริการ สามารถเรียกใช้ได้จากการเติมพารามิเตอร์ต่อท้าย URL พารามิเตอร์แรกก็คือเครื่องหมาย '?' จากนั้นพารามิเตอร์อื่นๆคือเครื่องหมาย '&' ตัวอย่างเช่น `.../martservice?type=attributes&dataset=hsapiens_gene_ensembl` และต้องระบุ *Virtual schema* ในฐานะข้อมูลนั้นๆด้วย (รูปที่ 2.25 ดาตาเซ็ทที่อยู่ภายในเส้นประคือ *Virtual schema*) การปรับแต่งเมตาดาตาของบริการไอโบบมาร์ทดังตารางที่ 2.2 และความหมายของพารามิเตอร์ใช้ในการปรับแต่งดังตารางที่ 2.3



รูปที่ 2.25 ตัวอย่างขอบเขตของ Virtual schema ใดๆของดาตามาร์ท [12]

ตารางที่ 2.2 เมตาดาตาและรายละเอียดการปรับแต่งของบริการไปโอมาร์ท

Metadata	Type	Description	virtualschema	mart	interface	dataset
Registry file	เอกสารรี ริสทรี	เอกสารรีริสทรีที่ใช้ ปรับแต่งการติดตั้ง มาร์ทเซอร์วิส	-	-	-	-
Datasets	ดาตาเซต	รายชื่อดาตาเซตที่พร้อม ให้บริการ	✓	✓	-	-
Dataset Configuration	การปรับแต่ง	รายละเอียดการ ปรับแต่งของดาตาเซต	✓	-	✓	-
Attributes	แอตทริบิวต์	ข้อมูลรายละเอียดของ แอตทริบิวต์ที่พร้อม บริการของดาตาเซต นั้นๆ	✓	-	✓	✓
Filters	ฟิลเตอร์	ฟิลเตอร์หรือเงื่อนไข การคิวรี	✓	-	-	✓

ตารางที่ 2.3 พารามิเตอร์ที่ใช้ในการปรับแต่งเมตาดาตา

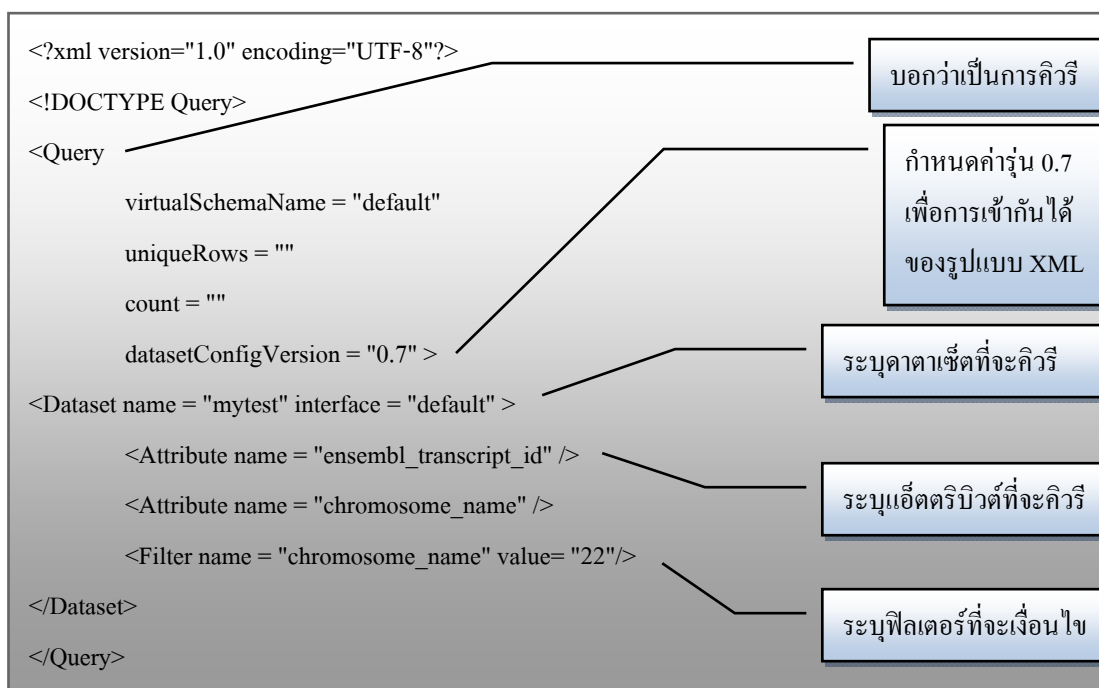
พารามิเตอร์	ความหมาย
virtualschema	ชื่อของโครงสร้างข้อมูลที่มีดาตาเซตนั้น (ค่าปริยายคือ 'default')
mart	ชื่อของดาตาเซตนั้นๆ
	ชื่อของผู้ใช้งานของดาตาเซต (ค่าปริยายคือ 'default')
interface	ชื่อของตัวเชื่อมประสานของดาตาเซตที่ใช้งาน (ค่าปริยายคือ 'default')
	ชื่อของผู้ใช้งานของดาตาเซต (ค่าปริยายคือ 'default')
dataset	ชื่อของผู้ใช้งานของดาตาเซต (ค่าปริยายคือ 'default')

2.2) การคิวรีข้อมูล (Queries)

คำสั่งคิวรีจะถูกส่งให้ระบบด้วยการร้องขอ *POST* โดยภายในโครงสร้างคำสั่งคิวรีจะมีพารามิเตอร์ที่ชื่อว่า *query* และค่าพารามิเตอร์ต่างๆ เป็นเอกสาร XML และผลลัพธ์ที่จะได้ก็คือที่ระบุไปเอกสาร XML นี้

- Query XML syntax

โครงสร้างเอกสาร XML ดังรูปที่ 2.26 คือตัวอย่างของคำสั่งการคิวรีซึ่งระบุต้องการ 1 ฟیلเตอร์ และ 2 แอ็ตทริบิวต์ ในการคิวรี 1 ดาตาเซ็ต



รูปที่ 2.26 ตัวอย่างของคำสั่งการคิวรีในดาตาเซ็ต

แท็ก *Query* ประกาศตัวเองว่าเป็นคิวรีข้อมูลแท็ก *datasetConfigVersion* กำหนดค่ารุ่น 0.7 เพื่อการเข้ากันได้ของรูปแบบ XML ที่กำหนดค่าการปรับแต่งต่างๆ ในดาตาเซ็ต แท็ก *Attribute* และ *Filter* สามารถกำหนดได้หลายค่าตามความต้องการของผู้ใช้ โดยแท็ก *Filter* สามารถกำหนดค่าได้ด้วย

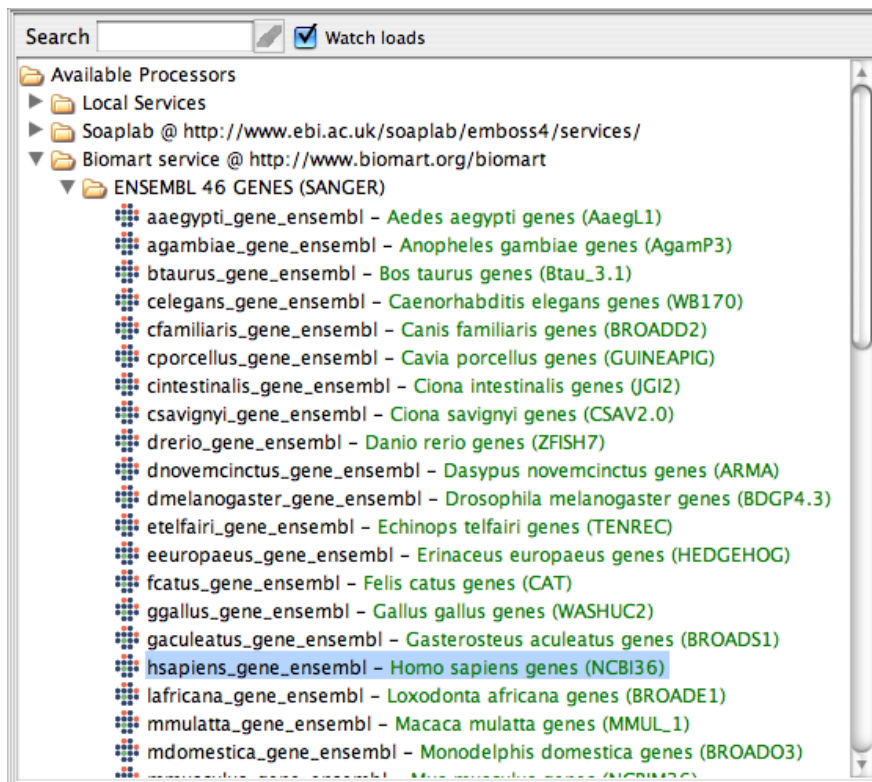
แอ็คตริบิวต์ที่เป็นผลลัพธ์จะถูกแสดงตามลำดับโครงสร้างข้างต้น โดยมีแท็บเป็น ตัวแยกโดยค่าปริยาย ถ้าแท็บ *uniqueRows* ถูกกำหนดค่าให้มีค่าเป็น 1 มาร์ทวิวกี่จะกำจัดข้อมูลที่ ซ้ำซ้อนกันให้ด้วย ขณะเดียวหากกำหนดค่าแท็บ *count* ให้มีเป็น 1 มาร์ทก็จะทำงานเพียงแค่นับ จำนวนแถวของผลลัพธ์เท่านั้น

2.6 การคิวรีข้อมูลด้วยบริการไบโอมาร์ทในโปรแกรมทาเวอร์น่า

ไบโอมาร์ทเป็นคลังข้อมูล (Data Warehouse) ที่รวบรวมข้อมูลด้านต่างๆทาง ชีววิทยาอย่างมากมายและซับซ้อน อนุญาตให้ผู้ที่สนใจเข้าถึงข้อมูลมากมายเหล่านั้นได้ เช่น บริการไบโอมาร์ทของ EBI มีข้อมูลจากฐานข้อมูล *Ensembl*, *VEGA*, *DbSNP*, *UniProt* และ *MSD* [12] และโปรแกรมทาเวอร์น่าเอง ก็มีกลไกการคิวรีของไบโอมาร์ทที่สามารถค้นหา และใช้งาน ข้อมูลเหล่านั้นได้เช่นเดียวกัน

2.6.1. การเพิ่มบริการไบโอมาร์ท

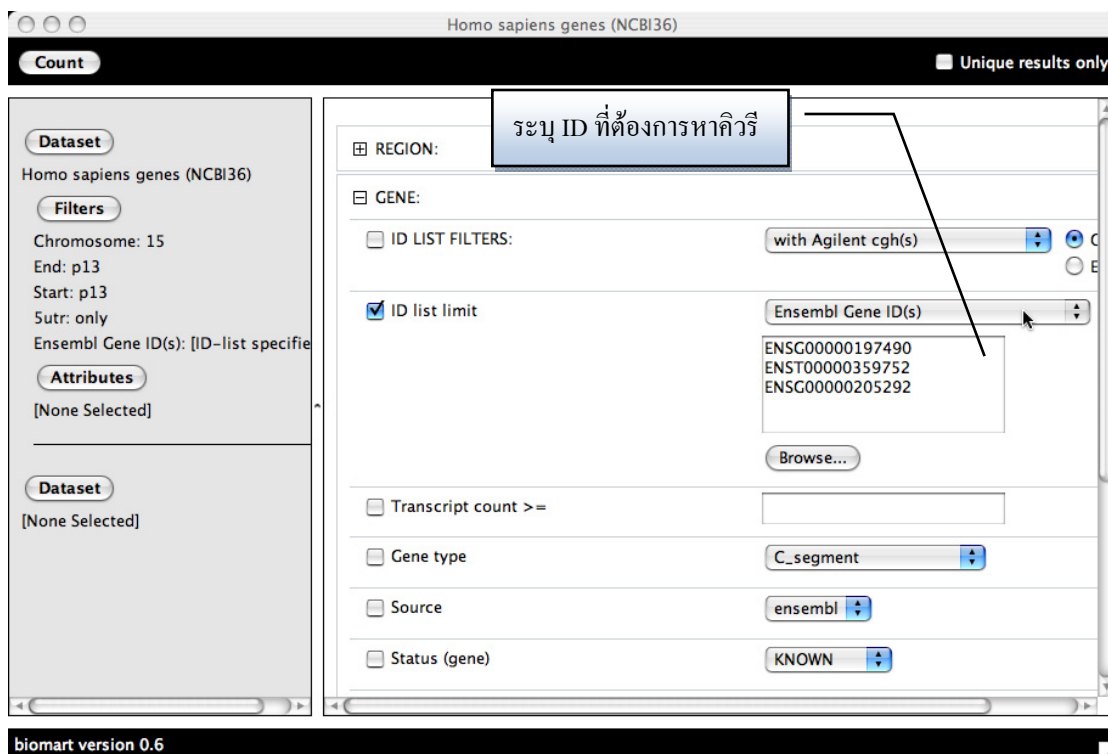
ผู้ใช้สามารถเพิ่มเอนพอยต์ของบริการไบโอมาร์ทใหม่ๆเข้าไปใน *Available service panel* ของโปรแกรมทาเวอร์น่า [8][9] หากเอนพอยต์สามารถเข้าถึงได้หรือมีการตอบกลับ (Responding) ก็จะได้แสดงรายชื่อโปรเซสเซอร์หรือบริการต่างๆใน *Available services panel* โดย สามารถแตกโครงสร้างต้นไม้เพื่อเลือกบริการที่ต้องใช้งานได้ ดังรูปที่ 2.27



รูปที่ 2.27 โครงสร้างต้นไม้ของบริการไบโอมาร์ทในโปรแกรมทาวเวอร์นา

- การปรับแต่งฟิลเตอร์ (Configuring Filters)

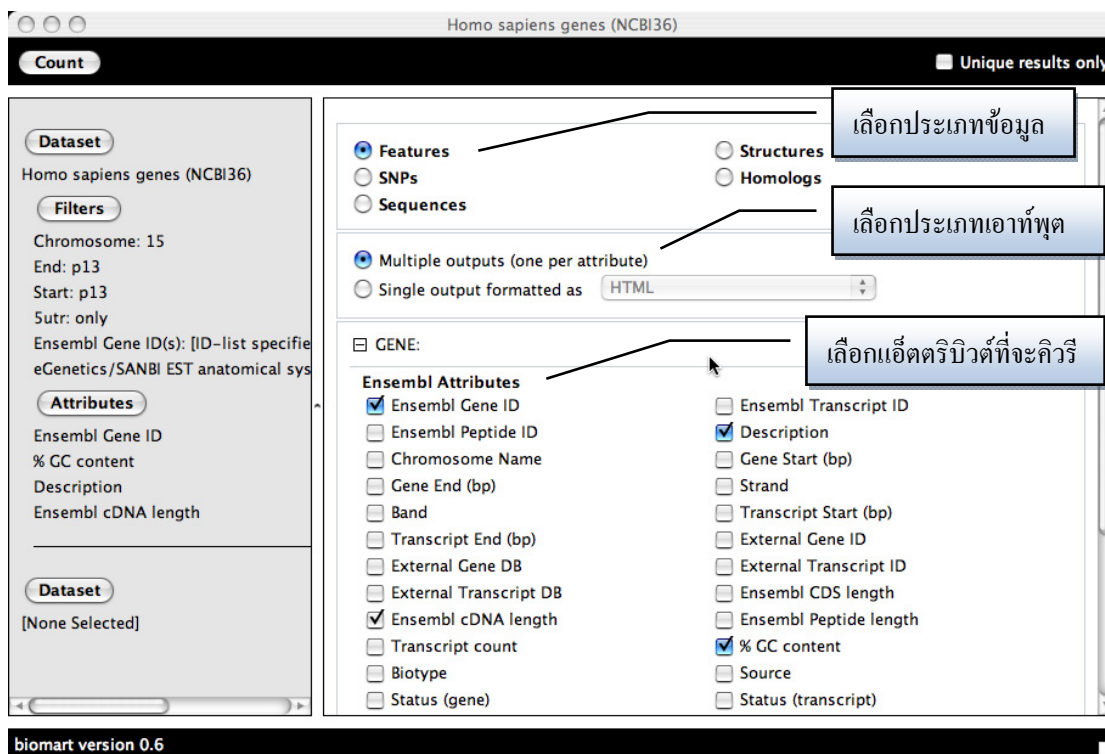
เมื่อเพิ่มบริการไบโอมาร์ทลงใน AME แล้ว สามารถปรับแต่งค่าที่ต่างๆให้ตรงตามความต้องการได้ผ่านหน้าต่าง 'Out of the box' ดังรูปที่ 2.28 โดยความเร็วในการทำงานของบริการไบโอมาร์ทจะขึ้นอยู่กับความเร็วของเครือข่ายที่ต่อออกสู่อินเทอร์เน็ต เนื่องจากหน้าต่างกราฟิก (GUI) ที่ทำหน้าที่ติดต่อกับผู้ใช้งานของไบโอมาร์ทนั้นต้องดึงข้อมูล (Fetch) จากเครื่องแม่ข่ายที่ให้บริการซึ่งผ่านเครือข่ายอินเทอร์เน็ต ซึ่งการปรับแต่งจะขึ้นอยู่กับความต้องการดึงหรือคิวรีข้อมูลที่ผู้ใช้งานสนใจ โดยฟิลเตอร์จะสามารถปรากฏให้เลือกได้ก็ต่อเมื่อมีการเลือกคาตาเซตแล้วเท่านั้น



รูปที่ 2.28 หน้าต่าง 'Out of the box' กำหนดค่าฟิลเตอร์ต่างๆ

2.6.2. การปรับแต่งแอตทริบิวต์ (Configuring Attributes)

แอตทริบิวต์หรือฟิลด์ข้อมูลที่ต้องการควิรีสามารถกำหนด และปรับแต่งที่หน้าจอของแอตทริบิวต์ ซึ่งจะแบ่งออกเป็นหน้าๆตามประเภทของข้อมูลที่สามารถควิรีได้ คือ *Feature*, *Structure*, *SNPs*, *Homologs* และ *Sequences* ดังแสดงในรูปที่ 2.29 การเลือกแอตทริบิวต์ สามารถเลือกได้เพียงประเภทเดียวเท่านั้นต่อหนึ่งดาตาเซต เนื่องจากแนวคิดการทำงานแบบอะตอมมิก หากผู้ใช้เลือกแอตทริบิวต์ของประเภทข้อมูลอื่นๆ แอตทริบิวต์ที่ถูกเลือกก่อนหน้าก็จะเอาออกโดยอัตโนมัติและแอตทริบิวต์ที่ได้เป็นผลลัพธ์ก็จะขึ้นอยู่กับเงื่อนไขของฟิลเตอร์นั้นๆ

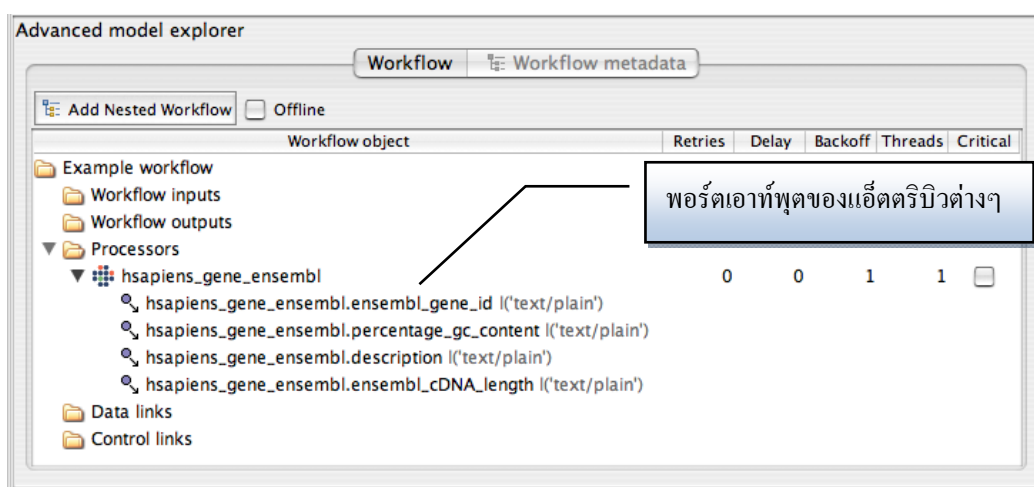


รูปที่ 2.29 หน้าต่าง 'Out of the box' กำหนดค่าแอตทริบิวต์

2.6.3. โหมดผลลัพธ์ของบริการไบโอมาร์ท (Result modes)

ผลลัพธ์ของการคิวรีด้วยบริการไบโอมาร์ทมีสองโหมดคือ แบบหลายเอาต์พุตและแบบเอาต์พุตเดียวซึ่งขึ้นอยู่กับ การเลือกจำนวนของแอตทริบิวต์ ซึ่งการเลือกจำนวนแอตทริบิวต์เพิ่มขึ้นหรือน้อยลง ก็จะทำให้จำนวนเอาต์พุตของบริการแปรผันตามไปด้วย

- หลายเอาต์พุต (Multiple Outputs) - แต่ละแอตทริบิวต์ที่ถูกเลือกจะกลายมาเป็นหนึ่งเอาต์พุตของบริการนั้นๆ ซึ่งจะมีชื่อเพียงพอที่สามารถคาดเดาได้ โดยไม่ยากนัก จากรูปที่ 2.30 แสดงบริการใน AME ที่ประกอบด้วยแอตทริบิวต์จำนวน 3 ตัว ซึ่งจะแสดงเป็น 3 เอาต์พุตพอร์ต
- เอาต์พุตเดียว (Single Output) - เช่นเดียวกับแบบหลายเอาต์พุต ขึ้นอยู่กับ การเลือกจำนวนแอตทริบิวต์ ในกรณีนี้คือเลือกเพียงหนึ่งแอตทริบิวต์ ก็จะได้เพียงหนึ่งเอาต์พุตปรากฏที่บริการนั้นๆ



รูปที่ 2.30 บริการไบโอมาร์ทที่มีหลายเอาต์พุต

2.7 สรุป

ในบทนี้ได้กล่าวถึงสถาปัตยกรรมของงานวิทยานิพนธ์ และหลักการสำคัญที่เกี่ยวข้องกับเวิร์คโฟลว์, โปรแกรมทาเวอร์น่า, เว็บเซอร์วิสและบริการต่างๆที่จะนำมาสร้างเป็นเวิร์คโฟลว์ รวมทั้งรายละเอียดเกี่ยวกับบริการ ไอ โบมาร์ทที่จะเป็นบริการหลักที่เลือกใช้ในวิทยานิพนธ์นี้

การค้นหบริการไบโอมาร์ทที่มีฟิลเตอร์และแอ็คทีวิตีใดที่ล้ำสมัย จะต้องค้นหาโดยวิธีการตรงไปตรงมา (Manual) คือต้องเรียกหน้าต่าง 'Out of the box' ของบริการไบโอมาร์ทขึ้นมาแล้วสังเกตว่าฟิลเตอร์และแอ็คทีวิตีใดที่เลือกไว้ก่อนหน้านั้น ยังปรากฏในหน้าต่าง 'Out of the box' หรือไม่ หากไม่มีปรากฏแล้ว ก็แสดงว่าฟิลเตอร์และแอ็คทีวิตีนั้นๆถูกผู้ให้บริการเอาออกหรือเปลี่ยนไปเป็นชื่อใหม่ เป็นต้น เหล่านี้ล้วนต้องอาศัยประสบการณ์ และการสังเกตจากการใช้บริการไบโอมาร์ทของผู้ใช้งานเอง ยิ่งกว่านั้น จำนวนของบริการไบโอมาร์ทที่อยู่ในเวิร์คโฟลว์, จำนวนแอ็คทีวิตีและฟิลเตอร์ที่เรียกใช้ในบริการไบโอมาร์ท, และความซับซ้อนของเวิร์คโฟลว์ ปัจจัยเหล่านี้ ล้วนมีผลต่อการตรวจสอบหรือการค้นหาสาเหตุที่ทำให้เวิร์คโฟลว์ทำงานได้ผลลัพธ์ที่ไม่ตรงกับความต้องการหรือไม่ให้ผลลัพธ์อะไรเลย

ดังนั้นขั้นตอนการตรวจสอบเวิร์คโฟลว์ก่อนการทำงานจึงมีความจำเป็น โดยจะมุ่งเน้นไปที่การตรวจสอบบริการไบโอมาร์ทในเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์ เพราะเป็นบริการหลักที่ใช้สำหรับงานวิจัยและเป็นกรณีศึกษาหลักของวิทยานิพนธ์นี้ ผู้วิจัยจึงได้เสนอ

แนวทางแก้ไขปัญหาโดยการพัฒนาเวิร์คโฟลว์ภายใต้สิ่งแวดล้อมของทาเวอร์นา ที่สามารถตรวจสอบทั้งฟิลเตอร์ และแอ็คติวิตีของบริการไปโอมาร์ทได้ เพื่อให้ข้อมูลที่มีนัยสำคัญต่อผู้ใช้งานจะได้มีความสะดวกและง่ายต่อการปรับแก้เวิร์คโฟลว์ ซึ่งจะเป็นการประหยัดเวลาในการดีบั๊กลงได้มากคงจะได้กล่าวในบทที่ 5 และบทที่ 6 ต่อไป

บทที่ 3

การเลือกบริการที่เหมาะสมและการสร้างบริการท้องถิ่นในการวิเคราะห์สนิปกึ่ง

3.1 บทนำ

โครงการมายคริดพัฒนาโปรแกรมทาเวอร์นา เพื่อให้ผู้ใช้สามารถเข้าถึงบริการหรือทรัพยากรอันหลากหลายต่างๆที่กระจุกกระจายกันอยู่บนอินเทอร์เน็ต โปรแกรมทาเวอร์นาสามารถเข้าถึงบริการได้กว่า 30,000 ตัวอย่างเช่น ชีวสารสนเทศ วิทยาศาสตร์สุขภาพ ชีวเคมี และดาราศาสตร์ เป็นต้น บริการจำนวนมากเหล่านี้ แต่ละบริการก็มีพารามิเตอร์และวิธีการเข้าถึงที่แตกต่างกันไป โครงการมายคริดออกแบบสถาปัตยกรรมการทำงานของโปรแกรมทาเวอร์นาให้ปกปิดความยุ่งยากซับซ้อนในการเข้าถึงทรัพยากรเหล่านี้ไว้ด้วยการทำงานแบบเวิร์คโฟลว์ ทำให้ผู้ใช้งานเข้าถึงได้โดยง่าย อีกทั้งบริการหรือทรัพยากรจะทำงานอยู่บนระบบที่มีฮาร์ดแวร์ที่ทรงประสิทธิภาพ เช่น ทำงานอยู่บนระบบคลัสเตอร์ขนาดหลายร้อยโหนด เป็นต้น

อย่างไรก็ตามงานวิจัยที่ใช้เทคโนโลยีเวิร์คโฟลว์ ยังจัดว่ามีลักษณะเฉพาะตัว จึงจำเป็นต้องพัฒนาเวิร์คโฟลว์ที่เหมาะสมกับงาน กล่าวคือ การเลือกบริการที่ทำงานได้ถูกต้องตรงกับความต้องการให้มากที่สุด เพราะไม่สามารถที่จะหาบริการที่ใช้ข้อมูลอินพุตและผลิตเอาต์พุตที่ต้องการสำหรับกระบวนการใดๆในงานวิจัยแบบเฉพาะจงจะมาใช้งานได้ทันที ซึ่งในบางกรณีอาจต้องมีการปรับแต่งรูปแบบของอินพุตและเอาต์พุตเหล่านั้นบ้างตามความจำเป็น นอกจากนี้ หากไม่สามารถหาบริการที่ต้องการได้ หรือบริการที่หาได้แล้วไม่สามารถตอบสนองความต้องการ เช่น การเกิดปัญหาการเกินเวลาของการส่งผ่านข้อมูลระหว่างบริการ ซึ่งเป็นปัญหาจากเครือข่ายหรือการติดต่อสื่อสารภายนอกโปรแกรมทาเวอร์นา และการไม่คงเส้นคงวาในการทำงานของเวิร์คโฟลว์ที่เรียกใช้บริการต่างๆที่อยู่บนอินเทอร์เน็ต (ทำงานได้บ้าง ไม่ได้บ้าง) เช่น เวิร์คโฟลว์ไม่ตอบสนองการส่งงานจากผู้ใช้ซึ่งเป็นปัญหาที่โปรแกรมทาเวอร์นาเอง

ดังนั้นบทนี้จึงนำเสนอการวาดและการเลือกบริการที่เหมาะสม โดยให้ข้อเสนอแนะที่เป็นรูปธรรมและวิธีการใช้เวิร์คโฟลว์ และเสนอบริการท้องถิ่นที่มีความจำเป็นใน

การทำงานสำหรับการแก้ปัญหาในการวิเคราะห์สลับคู่ [38] ที่เป็นงานวิจัยร่วมกับศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์เป็นกรณีศึกษา และอภิปรายผลการทดลอง

3.2 ขั้นตอนการเลือกบริการที่เหมาะสม

การเลือกบริการที่เหมาะสมตรงกับความต้องการในการสร้างเวิร์คโฟลว์ใดๆ ที่เสนอในวิทยานิพนธ์นี้ แบ่งวิธีการค้นหาและการเลือกบริการออกเป็น 2 วิธีคือ การค้นหาบริการจากเว็บพอร์ทัลของมายกริด และการค้นหาบริการจากอินเทอร์เน็ตโดยปลั๊กอินของโปรแกรมทาเวอร์น่า ดังนี้

3.2.1 การค้นหาบริการจากเว็บพอร์ทัลของมายกริด

การค้นหาบริการจากเว็บพอร์ทัลของมายกริด ซึ่งรวมไคลเรคทอรีของบริการหรือเว็บเซอร์วิสต่างๆที่โปรแกรมทาเวอร์น่าสามารถเข้าถึงได้ เรียกว่า *Biological Web Services* ดังรูปที่ 3.1 โดยมี URL คือ <http://www.mygrid.org.uk/wiki/Mygrid/BiologicalWebServices> ผู้ใช้สามารถเชื่อมต่อไปยังเว็บไซต์ของผู้ให้บริการ และสามารถทดลองใช้บริการเหล่านั้นกับเว็บเบราว์เซอร์ได้ทันที เช่น บริการ EBI หรือ ไบโอมาร์ท เป็นต้น

myGrid

You are here: myGrid wiki > MyGrid Web > TavernaWorkbench > BiologicalWebServices

Biological web services

This document lists the appropriate links of all known biological web services. Its purpose is to provide a registry of all bio services that can be accessed through a client such as [Taverna](#). It is organised by the provider or host of the service(s). To add your biological web services to the list, please contact [Franck Tanoh](#).

Many of the services have been [annotated](#) and can be retrieved and used from Taverna by accessing the [feta plugin](#) for Taverna.

Note that

- you don't have to list your services here before accessing them through Taverna.
- the services listed here are not by default in Taverna, you can manually add them to the Taverna workbench (see [Taverna Manual](#)).
- there are services available in many domains other than biology. This list only mentions those in biology.

Major data centres

- ↓ [EMBL-EBI, UK](#)
- ↓ [DDBJ, Japan](#)
- ↓ [NCBI, USA](#)
- ↓ [PDBJ, Japan](#)

Smaller projects and databases

- ↓ [Kanehisa Laboratory, Kyoto, Japan](#)
- ↓ [myGrid, Manchester, UK](#)
- ↓ [BASIS, University of Newcastle, UK](#)
- ↓ [BiMolecular Interaction Network Database, BIND, University of Toronto, Canada](#)
- ↓ [GeneCruiser, Broad Institute, Harvard-MIT, USA](#)

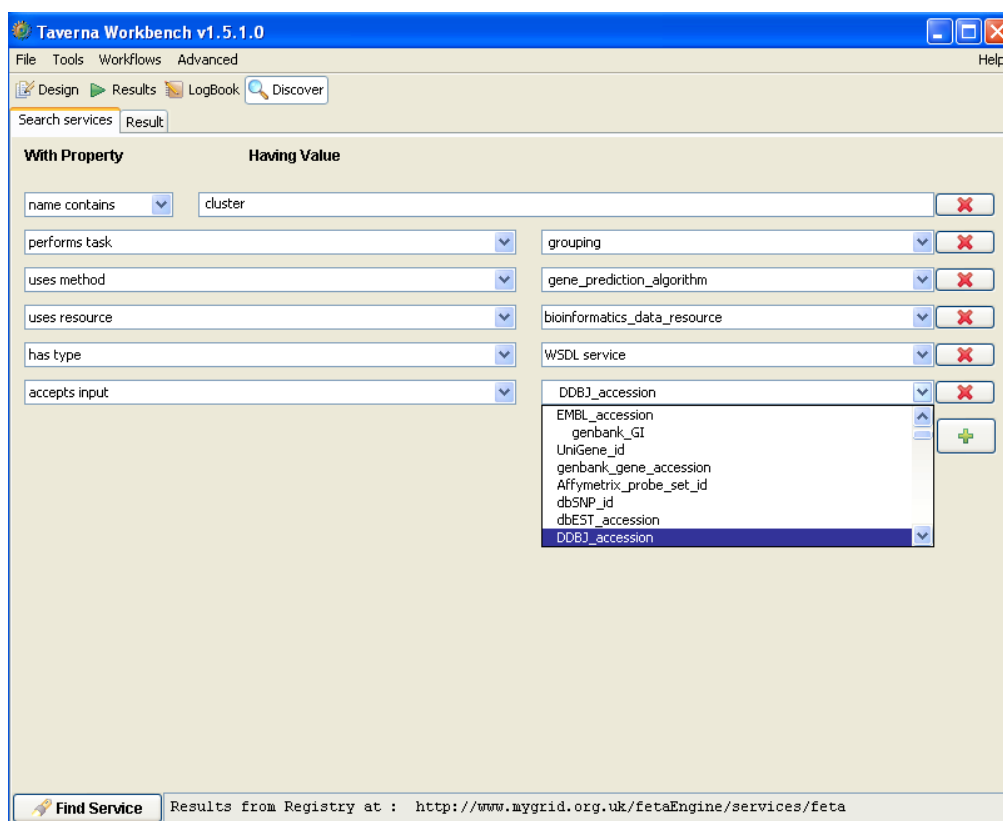
รูปที่ 3.1 เว็บไซต์ Biological web services ของมายกริด

3.2.2 การค้นหาบริการจากอินเทอร์เน็ตโดยปลั๊กอินของโปรแกรมทาวเวอร์นา

โปรแกรมทาวเวอร์นามีปลั๊กอินชื่อ *feta* เป็นตัวช่วยค้นหาบริการหรือเว็บเซอร์วิสต่างๆที่กระจัดกระจายบนเครือข่ายอินเทอร์เน็ต *feta* ค้นหาเว็บเซอร์วิสโดยใช้พื้นฐานจากคำอธิบายของหน้าที่ในการทำงานของบริการนั้นๆ (Description) ด้วย Service ontology ซึ่งจะเก็บอยู่ที่รีจิสทรีของบริการหรือ UDDI ของมายกริด หากพบบริการที่น่าจะเหมาะสมตรงตามความต้องการก็สามารถทดลองการใช้งานง่ายๆได้ ซึ่งอาจจะพบว่า ในงานหนึ่งๆนั้น มีบริการที่ทำงานคล้ายคลึงกันเป็นจำนวนมาก เนื่องจากการบริการเหล่านั้นได้รับการออกแบบโดยคำนึงถึงแนวคิดการทำงานแบบอะตอมมิกหรือไม่อาจจะแบ่งย่อยการทำงานลงได้อีก ซึ่งจะเป็นประโยชน์ต่อการดีบั๊กหรือตรวจสอบการทำงาน จากรูปที่ 3.2 เงื่อนไขในการค้นหาของ *feta* สามารถกำหนดได้หลายอย่างดังนี้ [8][9]

- ค้นหาบริการที่รับอินพุตในรูปแบบข้อมูลที่ต้องการ เช่น หมายเลขข้อมูลของฐานข้อมูล DDBJ

- ค้นหาบริการที่ผลิตเอาต์พุตในรูปแบบข้อมูลที่ต้องการ เช่น หมายเลขไอดีของสปีป
- ค้นหาบริการที่ทำหน้าที่ใดหน้าที่หนึ่งที่ต้องการ เช่น Grouping
- ค้นหาบริการที่มีอัลกอริทึมหรือ Method ที่ต้องการ เช่น Gene prediction algorithm หรืออัลกอริทึมในการทำนายยีน
- ค้นหาบริการที่มีฟังก์ชันการทำงานบางอย่างของโปรแกรมประยุกต์ เช่น Secondary protein structure prediction หรือการทำนายโครงสร้างชั้นที่สองของโปรตีน
- ค้นหาชนิดของบริการที่ต้องการ เช่น บริการ WSDL หรือบริการ Soaplab
- ค้นหาบริการจากชื่อและคำอธิบายการทำงานของบริการ เช่น Sapiens (มนุษย์)



รูปที่ 3.2 เงื่อนไขต่างๆในการค้นหาบริการของปลั๊กอิน Feta ในโปรแกรมทาวเวอร์น่า

จากเงื่อนไขต่างๆของการค้นหาโดย *Feta* จะเห็นว่า การระบุเงื่อนไขจะสามารถค้นหาบริการที่ตรงกับการทำงานของผู้ใช้ได้ เช่น ในวิทยานิพนธ์นี้ ระบุเงื่อนไขการค้นหาบริการชนิด Soaplab และระบุฟังก์ชันการทำงานว่า Translation หรือการถอดรหัสจากลำดับนิวคลีโอไทด์เป็นลำดับกรดอะมิโน เป็นต้น ในการใช้งานจริงนั้น ผู้ใช้อาจจะระบุชื่อหรือคำอธิบายการทำงานสั้นๆง่ายๆ และระบุประเภทของบริการ ก็เพียงพอที่จะหาบริการที่ตรงกับความต้องการจากรายการบริการต่างๆที่ *Feta* นำเสนอได้

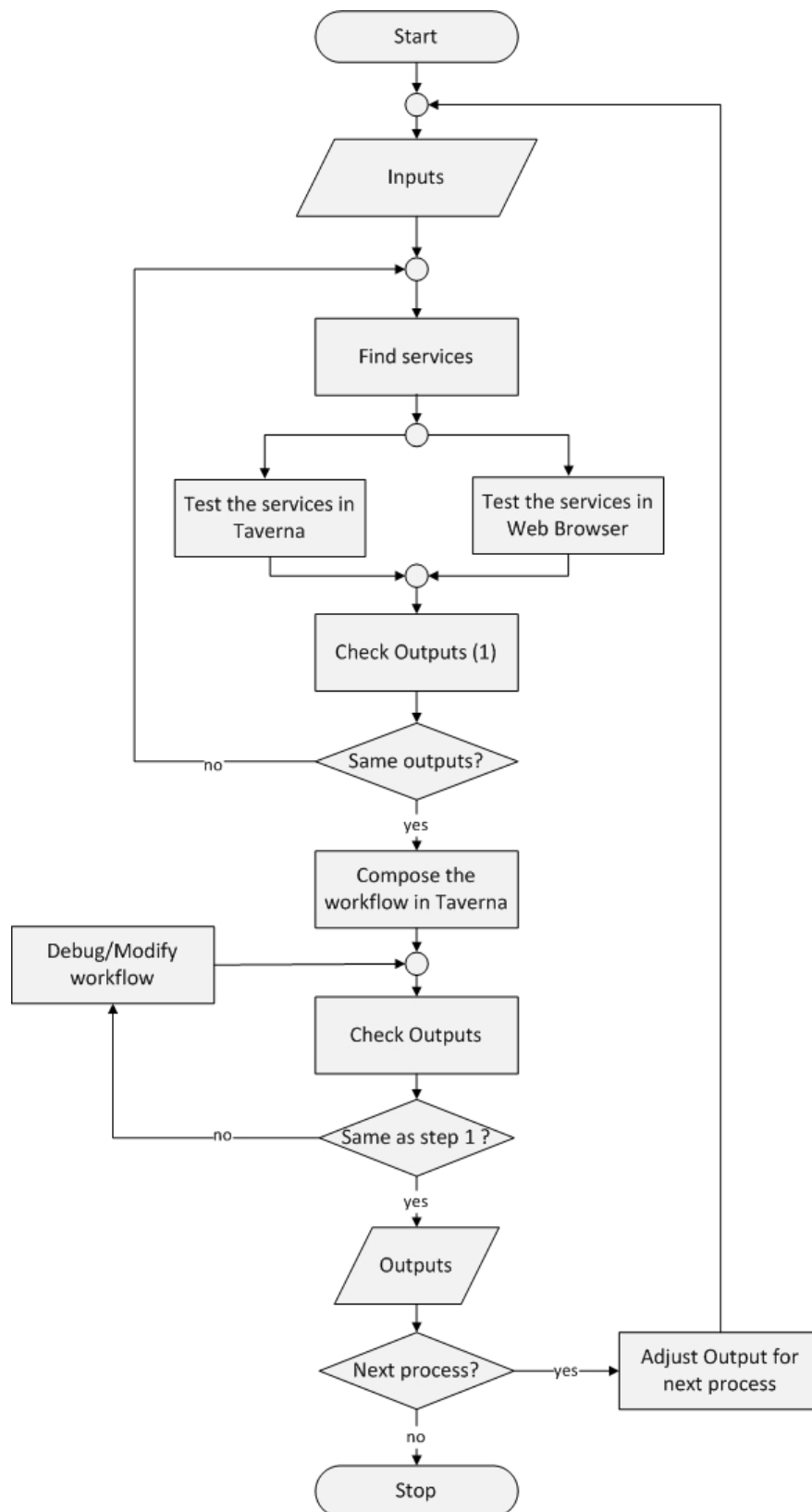
3.3 ขั้นตอนการพัฒนาเวิร์คโฟลว์

ขั้นตอนการพัฒนาเวิร์คโฟลว์ที่จะกล่าวในส่วนนี้ คือการออกแบบ Flowchart หรือ Data flow Diagram ของระบบงานวิจัยเสร็จสิ้นแล้ว ผู้วิจัยมองเห็นอินพุต การไหลของข้อมูลและเอาที่พุดโดยรวมจากการออกแบบได้ กระบวนการต่อไปคือ ขั้นตอนของการพัฒนาเวิร์คโฟลว์ของระบบงานตามที่ได้ออกแบบไว้แล้วดังรูปที่ 3.3

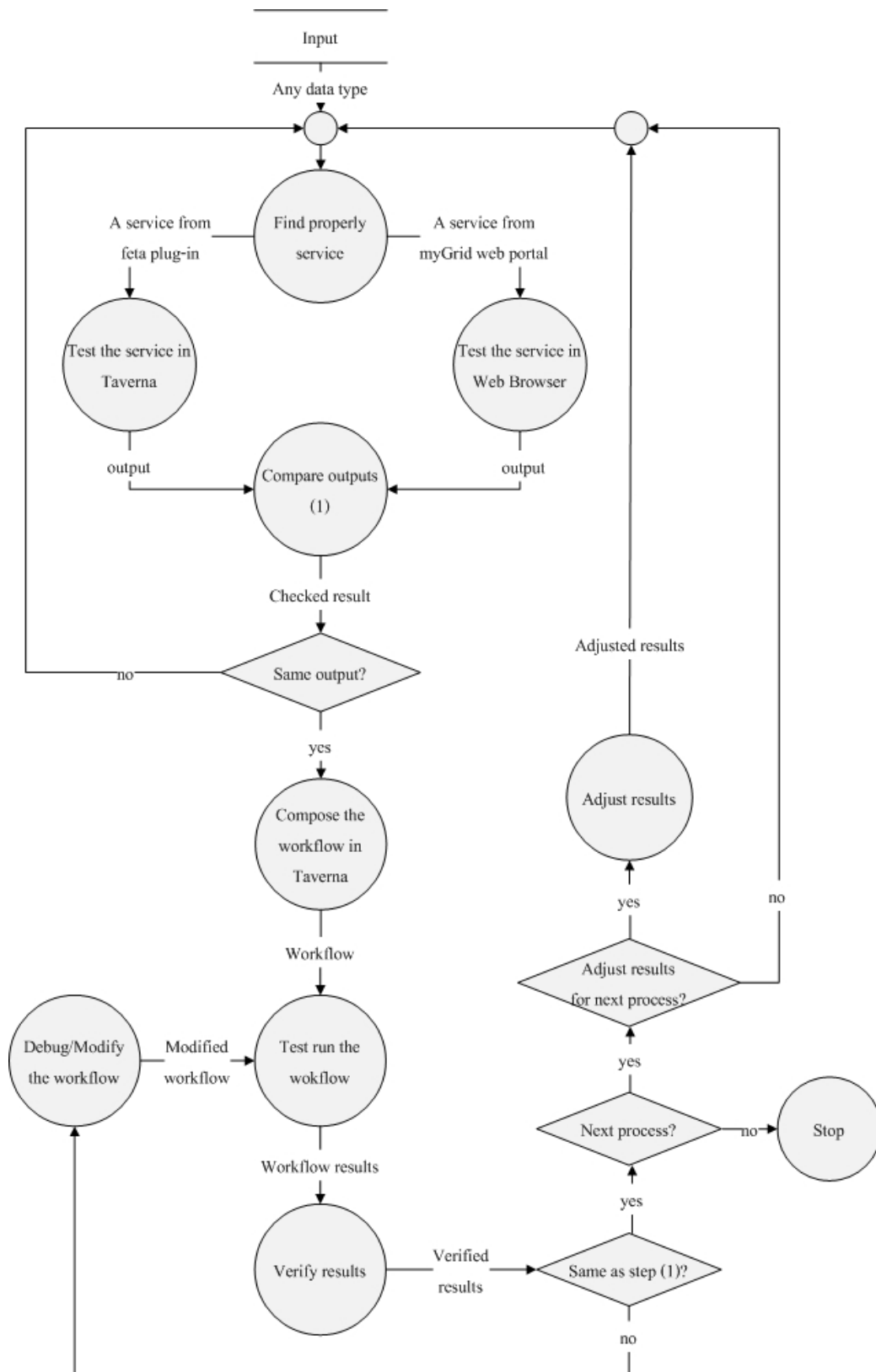
จากรูปที่ 3.3 สามารถอธิบายข้อมูลที่ไหลในกระบวนการพัฒนาเวิร์คโฟลว์ได้ด้วย Data flow diagram ดังรูปที่ 3.4 โดยขั้นเริ่มต้น มีข้อมูลอินพุตของงานวิจัยจากการออกแบบ Flowchart หรือ Data flow aiagram แล้ว จากนั้นเริ่มค้นหาบริการที่สามารถรองรับรูปแบบของอินพุตได้ หรืออาจจะปรับแต่งรูปแบบของอินพุตเพียงเล็กน้อย โดยไม่ใช้เวลาหรือแรงงานมากเกินไป ความจำเป็น ดังที่ได้กล่าวในหัวข้อที่ 3.2 ก็สามารถนำอินพุตที่ปรับแต่งนี้ไปใช้งานกับบริการที่เลือกมาได้

- ทดสอบการทำงานของบริการที่เลือกมานี้กับอินพุต ทั้งในโปรแกรมทาวเวอร์นา (ไม่ต้องวาดเป็นเวิร์ค โฟลว์ก่อน) และเว็บเบราว์เซอร์ แล้วนำผลลัพธ์การทำงานมาเปรียบเทียบกันว่า ถูกต้องตรงกันหรือไม่ เพื่อตรวจสอบความถูกต้องหรือดีบั๊กซึ่งกันและกัน และเป็นการเพิ่มความเชื่อมั่นในการทำงานของบริการที่นำมาประกอบเป็นเวิร์คโฟลว์ในขั้นตอนต่อไป
- หากผลการเปรียบเทียบไม่ถูกต้องตรงกัน แสดงว่าบริการที่เลือกมานั้นยังไม่อาจนำมาใช้งานได้ เพราะหากนำมาใช้แล้ว ในอนาคตบริการมีการเปลี่ยนแปลงก็จะทำให้ยากต่อการดีบั๊กหรือตรวจสอบ

- หากผลการเปรียบเทียบถูกต้องตรงกัน ขั้นตอนต่อไปก็นำบริการนี้ไปทดสอบสร้างเป็นเวิร์คโฟลว์ในโปรแกรมทาเวอร์น่า แล้วทดสอบการทำงาน และตรวจสอบผลลัพธ์การทำงานว่า ถูกต้องเหมือนกับที่ใช้งานกับเว็บเบราว์เซอร์หรือไม่
- หากผลลัพธ์การทำงานไม่ถูกต้อง แสดงว่ามีความผิดพลาดจากการสร้างเวิร์คโฟลว์หรือสาเหตุใดๆก็ได้ ผู้ใช้งานอาจจะตรวจสอบหรือปรับแต่ง Data flow ในเวิร์คโฟลว์อย่างละเอียดถี่ถ้วน ซึ่งโปรแกรมทาเวอร์เองก็มีเบรกพอยต์และรายงานสถานะของโปรเซสเป็นเครื่องมือดีบัก และลองทดสอบการทำงานซ้ำอีกครั้งหนึ่ง
- หากผลการทำงานจากเวิร์คโฟลว์ถูกต้องแล้ว ก็พิจารณาต่อว่าเอาท์พุตที่ได้มานี้ จะต้องไปเป็นอินพุตของกระบวนการหรือบริการถัดไปตามที่ได้ออกแบบไว้หรือไม่ ก็เริ่มต้นหาบริการที่เหมาะสมกับอินพุตนี้ต่อไป
- ทำตามกระบวนการดังกล่าวไปเรื่อยๆ จนกระทั่งครบทุกกระบวนการก็จะสามารถหาบริการที่ผลิตผลลัพธ์ที่ถูกต้องเหมาะสมกับงานได้



รูปที่ 3.3 ขั้นตอนการหาบริการที่เหมาะสมและการพัฒนาเวิร์คโฟลว์ที่เสนอในวิทยานิพนธ์นี้



รูปที่ 3.4 Data flow diagram อธิบายข้อมูลที่ไหลในกระบวนการพัฒนาเวิร์คโฟลว์

3.4 การออกแบบเวิร์คโฟลว์และเว็บเซอร์วิสสำหรับการวิเคราะห์สนิปกึ่ง

การสร้างเวิร์คโฟลว์เพื่อวิเคราะห์สนิปกึ่ง มีประโยชน์ในการพัฒนาสายพันธุ์ เป็นความร่วมมือกับสถานวิจัยจีโนมิกส์และชีวสารสนเทศ คณะวิทยาศาสตร์ มหาวิทยาลัยสงขลานครินทร์ ซึ่งเป็นกรณีศึกษาเบื้องต้นก่อนที่จะนำระบบไปปรับใช้งานกับกรณีศึกษาหลักต่อไป

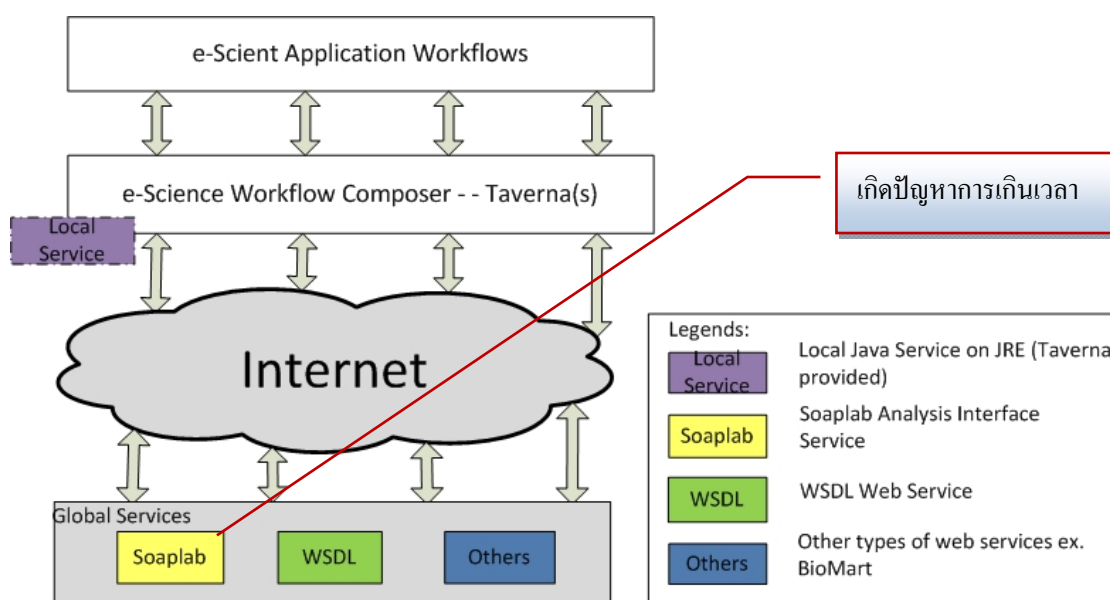
เวิร์คโฟลว์จะประกอบด้วยเว็บเซอร์วิสต่างๆ ที่ทำงานประสานกันไปตามหน้าที่ของเว็บเซอร์วิสนั้นๆ โดยการพัฒนาจะเริ่มตั้งแต่การค้นหาเว็บเซอร์วิสที่จะนำมาใช้งาน ซึ่งเลือกใช้ Emboss web service ของสถาบัน EBI [37] เป็นหลัก เนื่องจากโปรแกรมทาเวอร์น่าสามารถเข้าถึงได้อยู่แล้ว และมีการใช้งานที่ไม่ยุ่งยากซับซ้อน สามารถทดสอบความถูกต้องการทำงานของบริการ โดยการใช้เว็บเบราว์เซอร์ได้ ทำให้ง่ายต่อการตรวจสอบความถูกต้องของแต่ละเว็บเซอร์วิส ซึ่งเป็นไปตามกระบวนการที่นำเสนอในการเลือกบริการที่เหมาะสม

นอกจากนี้ ยังพัฒนาเว็บเซอร์วิสท้องถิ่น (Local web services) ในการทำงานด้านการสร้างโครงสร้างต้นไม้สายวิวัฒนาการ (Phylogenetics Tree) ขึ้นมาใช้งานเอง โดยใช้ Soaplab Analysis Tool เป็นเครื่องมือช่วยสร้างตัวเชื่อมประสานหรือ Web service interface ให้กับโปรแกรมประยุกต์หรือแอปพลิเคชันต่างๆ ด้านการสร้างโครงสร้างต้นไม้สายวิวัฒนาการเพื่อให้โปรแกรมทาเวอร์น่าเข้าถึงได้ เนื่องจากการใช้งานเว็บเซอร์วิสของสถาบัน EBI ในกระบวนการนี้ มีการวิเคราะห์และถ่ายโอนข้อมูลจำนวนมากในแต่ละบริการ ทำให้มีปัญหาคาการเกินเวลาดำเนินการและความไม่คงที่ในการทำงานส่งผลให้เวิร์คโฟลว์ทำงานได้ไม่สำเร็จ และจำเป็นต้องสร้างเว็บเซอร์วิสท้องถิ่นขึ้นมาใช้งานเอง ซึ่งในบทนี้จะแสดงรายละเอียดต่างๆ เกี่ยวกับกรณีศึกษาและเบื้องหลังของงาน วิธีการทำงานแบบเดิมและแบบใหม่ บริการหรือเว็บเซอร์วิสที่เลือกใช้งานในการพัฒนาเวิร์คโฟลว์ การแก้ไขปัญหาการเกินเวลาและความไม่คงเส้นคงวาในการทำงาน การพัฒนาบริการท้องถิ่นและผลการทดลอง ดังต่อไปนี้

3.4.1 กรณีศึกษาและเบื้องหลังของงาน

การวิจัยมุ่งเน้นการศึกษาชิ้นชื่อ *Amylase* ของกึ่งแซบวัย เพราะว่ามีรูปแบบที่ดีในการศึกษาด้านกระบวนการปรับตัว (Adaptation) และวิวัฒนาการ ตัวอย่างของ *Penaeus merguensis* (กึ่ง) เก็บจากแหล่งอาศัย 3 แห่งในอ่าวไทย คือ ทราย (TDE), สงขลา (SKE), และ

สุราษฎร์ธานี (SRE) [13] ซึ่งสื่อบริการจะเป็นสิ่งสำคัญต่อการศึกษาในอนาคตด้านการเติบโตและการทนต่อโรคของกุ้ง หลังจากที่หาบริการที่ต้องการใช้งานครบทุกบริการ ผู้วิจัยได้พัฒนาได้พัฒนาเวิร์คโฟลว์และแต่พบว่าเวิร์คโฟลว์ไม่สามารถทำงานได้สำเร็จเนื่องจาก Data flow ระหว่างเว็บเซอร์วิสเกิดการเกินเวลาในกระบวนการสร้างโครงสร้างต้นไม้อายวิวัฒนาการดังรูปที่ 3.5 ปัญหานี้เกิดเมื่อกระบวนการใช้ข้อมูลในการหาความน่าจะเป็นมากกว่า 2,000 ชุด ซึ่งเป็นปัญหาจากฝั่งผู้ให้บริการเว็บเซอร์วิสเอง ในที่นี้คือ EBİ แม้ว่าจะมีความเป็นไปได้ที่จะทำงานดังกล่าวแบบคัดลอกและวาง แต่เนื่องจากข้อมูลมีจำนวนมาก ทำให้ยากในการทำงานแบบคัดลอกและวางดังกล่าว อีกทั้งยังต้องใช้เวลาอีกด้วย



รูปที่ 3.5 ภาพรวมการติดต่อสื่อสารของโปรแกรมทาเวอร์นาและจุดที่เกิดปัญหาการเกินเวลา

วิธีการแก้ไขปัญหาวินิจฉัยที่วิทยานิพนธ์นี้นำเสนอคือ การสร้างบริการท้องถิ่นที่สามารถติดต่อสื่อสารกับเวิร์คโฟลว์ได้ โดยจะเน้นไปที่กระบวนการสร้างโครงสร้างต้นไม้อายวิวัฒนาการของลำดับนิวคลีโอไทด์ (Nucleotide) และกรดอะมิโน (Amino Acid) ซึ่งตัวช่วยที่สามารถติดต่อสื่อสารนี้เรียกว่า Web service interfaces โดยจะทำหน้าที่เชื่อมประสานให้กับแต่ละกระบวนการของงานในการสร้างโครงสร้างต้นไม้อายวิวัฒนาการในเวิร์คโฟลว์ ซึ่งสามารถเข้าถึงได้โดยโปรแกรมทาเวอร์นา ตัวเชื่อมประสานเว็บเซอร์วิสในวิทยานิพนธ์ใช้โปรแกรม Soaplab Analysis Tool เป็นเครื่องมือในการพัฒนา

เวิร์คโฟลว์ที่พัฒนาในวิทยานิพนธ์ประกอบด้วยระบบงานดังนี้คือ BLAST (Basic Local Alignment and Search Tool) [41], การเรียงลำดับแบบหลายลำดับของลำดับนิวคลีโอไทด์ และกรดอะมิโน (Multiple sequence alignment), การทำนายโครงสร้างชั้นที่สองของโปรตีน (Prediction of the protein secondary structure) และการทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์และกรดอะมิโน (Prediction of the phylogenetic tree) ผู้วิจัยใช้โปรแกรมทาเวอร์น่าในการรวบรวมเว็บเซอร์วิส ที่เกี่ยวข้องกับระบบงานเข้ามาวาดเป็นเวิร์คโฟลว์ได้แก่ บริการ Web API for Biology จาก DDBJ [40], บริการ EBI [37], บริการท้องถิ่น Java และบริการท้องถิ่นที่สร้างขึ้นมาเอง ความถูกต้องของผลการทำงานในแต่ละบริการตรวจสอบโดยศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์

3.4.2 วิธีการทำงานแบบเดิม (Copying and Pasting)

การวิเคราะห์สลับไปในวิธีการทำงานแบบเดิมหรือการคัดลอกแล้ววาง (Copying and Pasting) เป็นขั้นตอนการทำงานด้วยมือโดยการคัดลอกและวางสลับกันระหว่างแอปพลิเคชันทางด้านชีวสารสนเทศศาสตร์ ซึ่งทำงานแบบออฟไลน์กับเว็บแอปพลิเคชันที่ทำงานบนอินเทอร์เน็ต การวิเคราะห์สลับแบบดั้งเดิมนี้อาศัยแอปพลิเคชัน และเว็บแอปพลิเคชันออนไลน์ต่างๆ ดังต่อไปนี้ [14]

1. โปรแกรม BLAST กับลำดับนิวคลีโอไทด์ เป็นการค้นหาข้อมูลของลำดับนิวคลีโอไทด์ของกิ้งที่ได้จากสกัดจากห้องปฏิบัติการด้วยโปรแกรม BLAST ที่เว็บไซต์ของ NCBI เพื่อหาว่าลำดับนิวคลีโอไทด์ที่สนใจเป็นยีน *Amylase* หรือไม่
2. โปรแกรม Clustal X และ GENEDOC กับลำดับนิวคลีโอไทด์ เป็นการทำ Multiple alignment หรือการเรียงลำดับแบบหลายลำดับของลำดับนิวคลีโอไทด์โดยใช้โปรแกรม Clustalx และ Genedoc ที่เครื่องคอมพิวเตอร์ท้องถิ่นส่วนบุคคลเพื่อดูว่าลำดับต่างๆ ที่นำมาวิเคราะห์มีส่วนคล้ายคลึงกันมากน้อยเพียงใด
3. โปรแกรม Protein sequence analysis กับลำดับนิวคลีโอไทด์ เป็นการ Translation หรือถอดรหัสจากลำดับนิวคลีโอไทด์ซึ่งเป็น DNA ไปสู่ลำดับกรดอะมิโนโดยใช้โปรแกรม Protein Sequence Analysis (PSA) ที่เครื่องคอมพิวเตอร์ท้องถิ่นส่วนบุคคลแล้วนำลำดับอะมิโนที่ได้ไปสู่กระบวนการวิเคราะห์ต่อไป

4. โปรแกรม Clustal X และ GENEDOC กับลำดับกรดอะมิโน เป็นการทำให้ Multiple alignment การเรียงลำดับแบบหลายลำดับของลำดับกรดอะมิโนที่ได้จากขั้นตอนที่ 3 โดยใช้โปรแกรม Clustalx และ Genedoc ที่เครื่องคอมพิวเตอร์ท้องถิ่นส่วนบุคคลเพื่อดูว่าลำดับต่างๆที่นำมาวิเคราะห์มีส่วนคล้ายคลึงกันมากน้อยเพียงใด

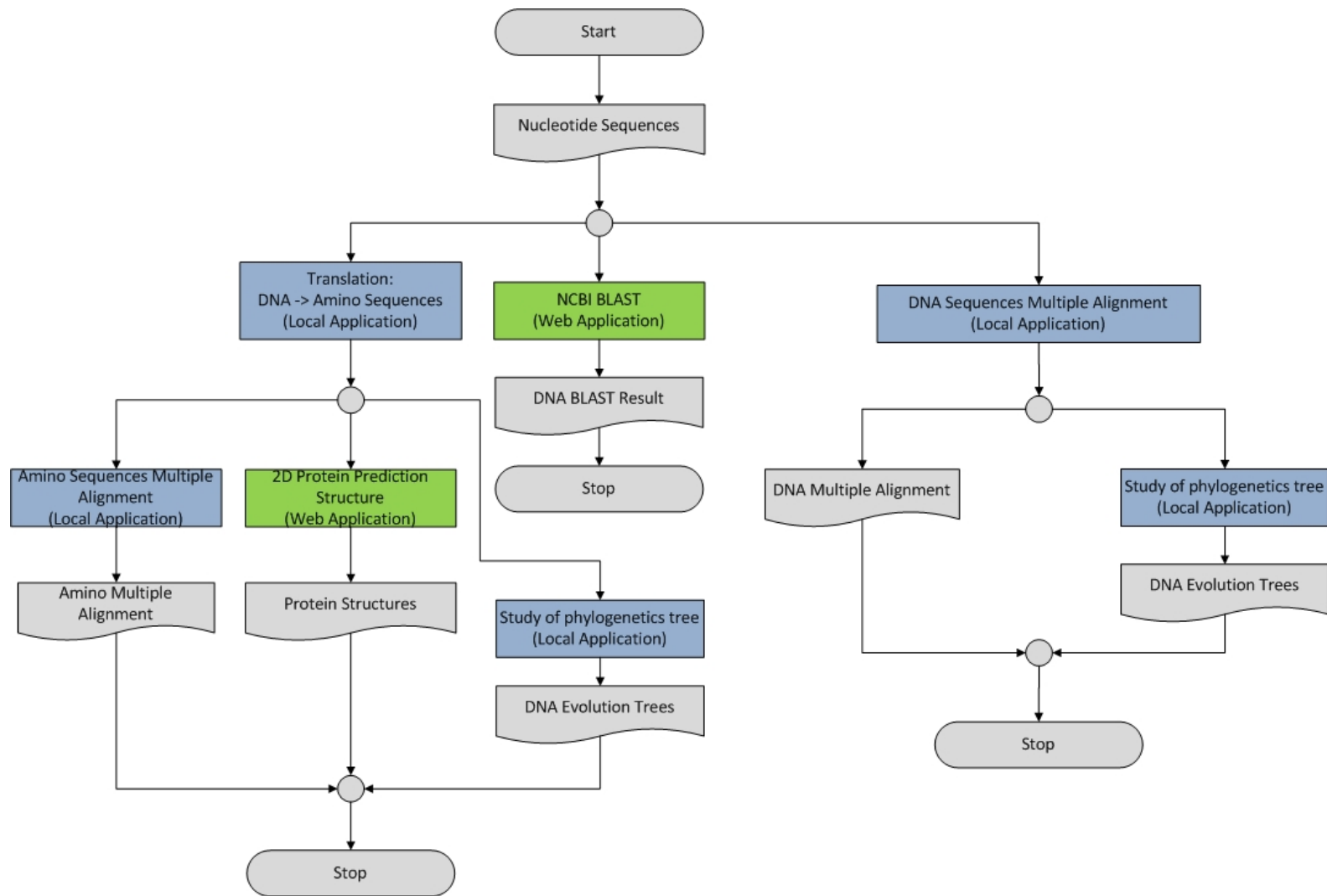
5. โปรแกรม BioEdit กับลำดับนิวคลีโอไทด์และลำดับกรดอะมิโน เป็นการศึกษาของความสัมพันธ์ทางวิวัฒนาการของสายพันธุ์ โดยวิเคราะห์จากโครงสร้างต้นไม้หรือ Phylogenetics tree ซึ่งสร้างได้จากโปรแกรม BioEdit ที่เครื่องคอมพิวเตอร์ท้องถิ่นส่วนบุคคลของนักวิจัยเอง โดยจะวิเคราะห์ทั้งลำดับนิวคลีโอไทด์และลำดับกรดอะมิโน

6. เว็บแอฟพลิเคชันชื่อโปรแกรม PSIPRED ที่ <http://bioinf.cs.ucl.ac.uk/psipred/> กับลำดับกรดอะมิโน เป็นการทำนายโครงสร้างขั้นที่สองของโปรตีน (Protein secondary structure) โดยใช้ข้อมูลจากลำดับกรดอะมิโนซึ่งมีเว็บแอฟพลิเคชันที่ให้บริการคือ PSIRRED

จากขั้นตอนของการทำงานทั้ง 1-6 ขั้นตอน สามารถสรุปความสัมพันธ์ของแต่ละกระบวนการได้รูปที่ 3.6 ซึ่งแสดงการใช้งานแอฟพลิเคชันบนเว็บเพจด้วยบล็อกสีเขียว และการใช้งานแอฟพลิเคชันที่เครื่องคอมพิวเตอร์ท้องถิ่นส่วนบุคคลด้วยบล็อกสีน้ำเงิน

3.4.3 บริการหรือเว็บเซอร์วิสที่เลือกใช้งานและการพัฒนาเวิร์กโฟลว์

รูปที่ 3.6 แสดงภาพรวมของการทำงาน และเห็นลักษณะการไหลของข้อมูลในกระบวนการวิเคราะห์สลับ โดยข้อมูลที่เป็นอินพุตให้กับระบบก็คือลำดับนิวคลีโอไทด์ของกุ้งที่ได้จากห้องปฏิบัติการ ในขั้นตอนนี้สามารถเริ่มหาบริการที่เหมาะสมกับความต้องการมาใช้งานตามกระบวนการในรูปที่ 3.3 และพัฒนาเวิร์กโฟลว์ด้วยโปรแกรมทาวเวอร์นา ในการพัฒนาเวิร์กโฟลว์นี้ได้ใช้บริการท้องถิ่นของ Java ในการจัดการ Data flow ระหว่างบริการต่างๆด้วย ในที่นี้คือ 'Split string into string list' (บล็อกสีม่วงในรูปที่ 3.7) ทำหน้าที่ส่งลำดับเบสไปเข้าไปทำงานกับโปรแกรม BLAST ที่ละลำดับดังแสดงในตารางที่ 3.1



รูปที่ 3.6 ลำดับของกระบวนการวิเคราะห์ขั้นต้นของกิ้ง

ตารางที่ 3.1 แสดงตัวอย่างการทำงานของบริการ ‘Split string into string list’

Input Type	String
Input Example	‘ATCGATCGATCGATCGATCG\ ATCGATCGATCGATCGATCG\ ATCGATCGATCGATCGATCG\ ’
Output Type	String list
Output Example	‘ATCGATCGATCGATCGATCG’ ‘ATCGATCGATCGATCGATCG’ ‘ATCGATCGATCGATCGATCG’
Expression	‘\n’

บริการหรือเว็บเซอร์วิสหลักที่เลือกมาใช้งานเป็นบริการของชุดโปรแกรม Emboss จากสถาบัน EBI โดยมีบริการที่ทำงานได้ตรงตามความต้องการและเป็นที่รู้จักกันอย่างแพร่หลาย ตารางที่ 3.2 แสดงรายชื่อบริการของชุดโปรแกรม Emboss ที่เกี่ยวข้องและตรงตามความต้องการของกระบวนการทำงาน (บล็อกสี่เหลี่ยมในรูปที่ 3.7) กระบวนการทำงานของโปรแกรม BLAST เรียกใช้บริการจาก DDBJ (บล็อกสี่เหลี่ยมในรูปที่ 3.7) ซึ่งสามารถระบุได้ว่า จะค้นหาลำดับนิวคลีโอไทด์ในฐานข้อมูลอะไรบ้าง ในงานวิจัยนี้เลือกใช้ฐานข้อมูล DDBJNEW [14] โปรแกรม BLAST ของ DDBJ สามารถเข้าถึงได้ด้วยโปรแกรมทาวเวอร์นาและเว็บเบราว์เซอร์ นอกจากนี้ยังสามารถทดลองการใช้งานจากเว็บเบราว์เซอร์ได้เช่นเดียวกันกับบริการของ EBI เวอร์คโพล์ที่พัฒนาเสร็จแล้ว แสดงดังรูปที่ 3.7 โดยมีพารามิเตอร์ของแต่ละบริการหรือเว็บเซอร์วิสที่กำหนดไว้ในเวอร์คโพล์ดังตารางที่ 3.3

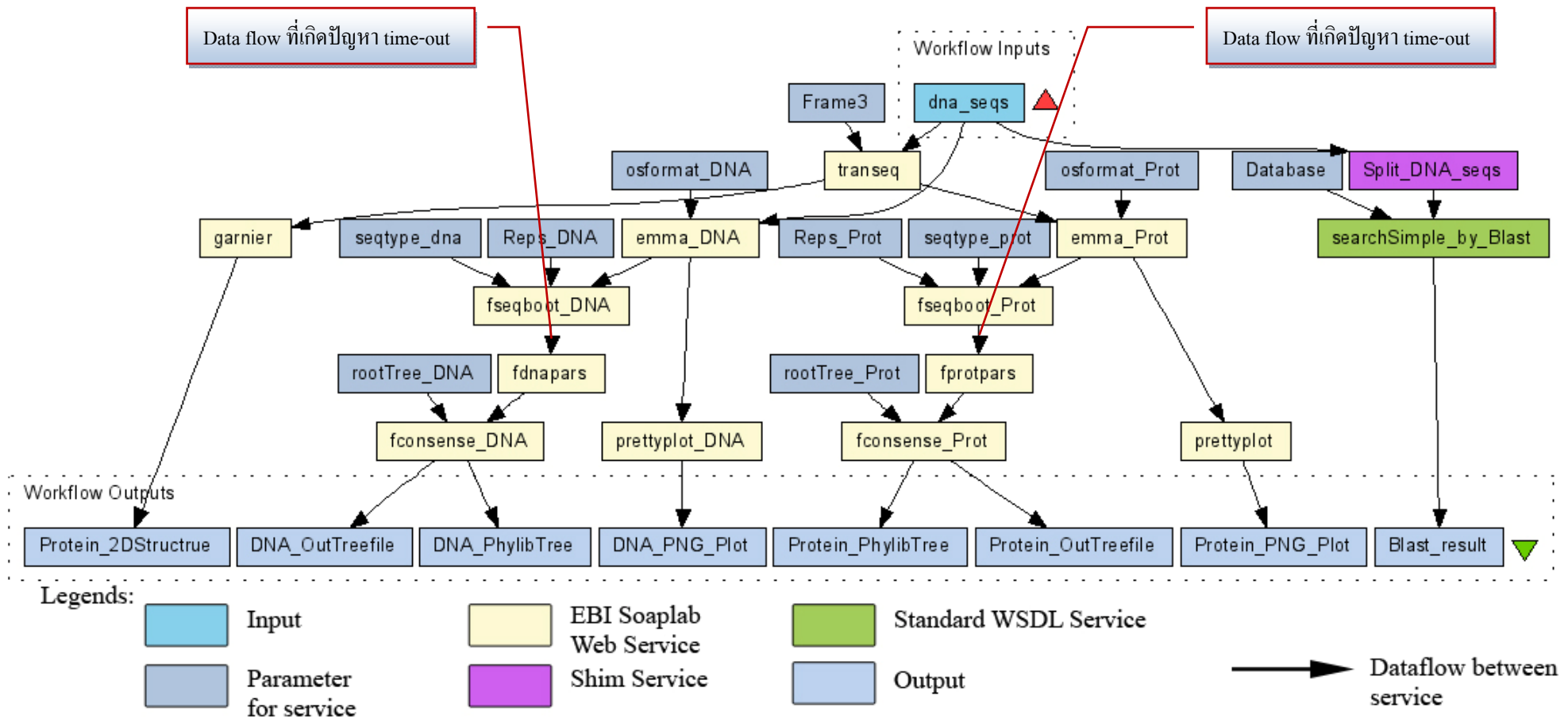
ในกระบวนการทำนาย หรือการสร้างโครงสร้างต้นไม้สายวิวัฒนาการของสิ่งมีชีวิตใดๆ (ในวิทยานิพนธ์นี้คือกุ้งแชบ๊วย) จะใช้บริการ *fseqboot* หรืออัลกอริทึมในการทำ Bootstrapping replicate กับลำดับข้อมูล *fseqboot* จะอ่านชุดข้อมูลและสร้างชุดข้อมูลของลำดับนิวคลีโอไทด์หรือกรดอะมิโนหลายๆชุดตามที่ผู้ใช้งานระบุโดยวิธีการแบบสุ่ม [37]

ตารางที่ 3.2 บริการ EBI ที่ใช้ในการพัฒนาเวิร์กโฟลว์วิเคราะห์ของกุ่ม

ชื่อบริการ EBI	หน้าที่การทำงาน
transeq	ถอดรหัสจากลำดับนิวคลีโอไทด์ (ดีเอ็นเอ) เป็นลำดับกรดอะมิโน (โปรตีน)
garnier	ทำนายโครงสร้างชั้นที่สองของโปรตีน
emma_[DNA/Prot]	หาความคล้ายกันของลำดับ (Multiply alignment) ด้วยโปรแกรมประยุกต์ ClustalW
prettyplot	แสดงผลการเรียงของลำดับด้วยสีและเครื่องหมายบ่งความเหมือนและแตกต่าง
fseqboot_[DNA/Prot]	อัลกอริทึมในการทำ Bootstrapping replication
fdnapars	อัลกอริทึม DNA parsimony เพื่อเตรียมข้อมูลสำหรับการวิเคราะห์โครงสร้างต้นไม้สายวิวัฒนาการ
fprotpars	อัลกอริทึม Protein parsimony เพื่อเตรียมข้อมูลสำหรับการวิเคราะห์โครงสร้างต้นไม้สายวิวัฒนาการ
fconsense_[DNA/Prot]	หาโครงสร้างต้นไม้สายวิวัฒนาการที่มีความเป็นไปได้มากที่สุด (Majority-rule and strict consensus tree)

ตารางที่ 3.3 พารามิเตอร์ของแต่ละบริการหรือเว็บเซอร์วิสที่กำหนดในเวิร์คโฟลว์

พารามิเตอร์	คำอธิบาย	ค่าที่กำหนด
Database	ชื่อของฐานข้อมูลที่จะให้โปรแกรม BLAST เข้าไปค้นหาลำดับเบส	DDBJNEW
Frame3	หมายเลขเฟรม 3 สำหรับการถอดรหัส [13][14]	3
osformat_[DNA/Prot]	รูปแบบลำดับของเอาต์พุตสำหรับวิเคราะห์โครงสร้างต้นไม้สายวิวัฒนาการ	phylip
seqtype_[DNA/Prot]	รูปแบบของการทำ Bootstrapping replicate ของลำดับเบส	ดีเอ็นเอ = d โปรตีน = p
Reps_[DNA/Prot]	จำนวนครั้งของ Bootstrapping replicate ในสร้างข้อมูลซ้ำๆ	ผันแปรตามกรณีทดสอบ
rootTree_[DNA_Prot]	ให้สร้างโครงสร้างต้นไม้สายวิวัฒนาการแบบมีราก	Yes



รูปที่ 3.7 เวิร์กโฟลว์แรกๆที่พัฒนาขึ้นโดยใช้บริการต่างๆที่กระจายอยู่บนอินเทอร์เน็ต

อินพุตของเวิร์คโฟลว์คือลำดับนิวคลีโอไทด์ของกิ่งแซบวัยในรูปแบบ Fasta โดยได้จากห้องปฏิบัติการ ซึ่งมีจำนวน 17-25 ลำดับ แต่ลำดับประกอบด้วย 250 คู่เบสที่ได้มาจากส่วน Exon7-10 ของโครโมโซมกิ่ง ตัวอย่างอินพุตลำดับนิวคลีโอไทด์รูปแบบ Fasta ดังรูปที่ 3.8

```
>SRE26
GTGGCGAAGCCATATCCAGCGGGCAGTATGTTGGCAACGGTCGTGTGACGGAGTTCAGGTACGGCAAGTAC
CTGGGCGAGGCCTCCGCGGCAACAACCAGCTGAAATACCTCAACAACCTTCGGCGAAGGTTGGGGCATGAT
TGACCGGCATGACGCCCTGGTCTTCATTGACAACCACGACAACCAGAGAGGGCCATGGTGCTGGAGGAGACA
TGATCCTTACTTTCCGTGTCTCTAAGTGGTACAAGGA
>TDE8
GTGGCGAAGCCATATCCAGCGGGCAGTATGTTGGCAACGGTCGTGTGACGGAGTTCAGGTACGGCAAGTAC
CTGGGCGAGGCCTCCGCGGCAACAACCAGCTGAAATACCTCAACAACCTTCGGCGAAGGTTGGGGCATGAT
TGACCGGCATGACGCACTGGTCTTCATTGACAACCACGACAACCAGAGAGGGCCATGGTGCTGGAGGAGACA
TGACCTTACTTTCCGTGTCTCTAAGTGGTACAAGGA
>SKE42
GTGGCGAAGCCATATCCAGCGGGCAGTATGTTGGCAACGGTCGTGTGACGGAGTTCAGGTACGGCAAGTAC
CTGGGCGAGGCCTCCGCGGCAACAACCAGCTGAAATACCTCAACAACCTTCGGCGAAGGTTGGGGCATGAT
TGACCGGCATGACGCACTGGTCTTCATTGACAACCACGACAACCAGAGAGGGCCATGGTGCTGGAGGAGACA
TGATCCTTACTTTCCGTGTCTCTAAGTGGTACAAGGA
```

รูปที่ 3.8 ตัวอย่างอินพุตลำดับนิวคลีโอไทด์รูปแบบ Fasta ของกิ่งแซบวัย

กระบวนการสร้างหรือทำนายโครงสร้างต้นไม้สายวิวัฒนาการ จะเรียกบริการ *fseqboot* สำหรับการทำให้ Bootstrapping replicate, *fdnapars* และ *fprotpars* สำหรับการหาของ อัลกอริทึม DNA parsimony และ Protein parsimony ตามลำดับและ *fconsense* สำหรับการหา โครงสร้างต้นไม้สายวิวัฒนาการที่มีความเป็นไปได้มากที่สุด

ตารางที่ 3.4 แสดงผลการทดสอบของเวิร์คโฟลว์ด้วยจำนวนของการทำซ้ำ (Number of bootstrapping replicates) ที่จำนวนต่างๆคือ 100, 200, 500, 1,000 และ 2,000 ตามลำดับ และตรวจวัดเวลาในการทำงาน พบว่าสำหรับการกำหนด Bootstrapping replicates ที่จำนวน 2,000 ชุด ทำให้ Data flow ระหว่างบริการ *fseqboot* -> *fdnapars* และ *fseqboot* -> *fprotpars* ไม่สามารถทำงานได้เนื่องจากเกิดเกินเวลาหรือเวิร์คโฟลว์ไม่ตอบสนองในการทำงานใดๆ (Not responding) เกิดความไม่คงเส้นคงวาในการทำงานช่วงส่งถ่ายข้อมูลเอาท์พุตระหว่างกันของเว็บเซอร์วิสดังกล่าว

ข้อความแสดงความผิดพลาดดังกล่าวที่โปรแกรมทาวเวอร์นารายงานในรูปที่ 3.9 หมายถึง มีปัญหาเกี่ยวกับการสื่อสารภายในบริการ Sopalab ของ EBI อาจเนื่องด้วยใช้เวลาในการ ค้นหาข้อมูลเกินค่าเวลาการทำงาน โดยปริยาย (Default time-out) ภายในของบริการ EBI ที่ตั้งค่าไว้ ซึ่งเป็นปัญหาภายในของ EBI เอง เราไม่สามารถเข้าไปแก้ไขได้ จึงจะต้องหาแนวทางแก้ไขที่ฝั่งของเราเอง

ตารางที่ 3.4 กรณีทดสอบและผลการทดสอบ

จำนวนชุดข้อมูลของการทำซ้ำ (No. of bootstrapping replicates)	เวลาที่ใช้ทำงาน (นาที)
100	02.50
250	12.20
500	20.29
1,000	20.40
2,000	ล้นเหลว (เกินเวลา)

```
ServiceError Message="Task failed due to problem invoking soaplab service"
TimeStamp="Feb 6, 2008 11:27:19 AM"
org.embl.ebi.SooplabShare.SooplabException:
Internal communication failed. (in fetchResults)
```

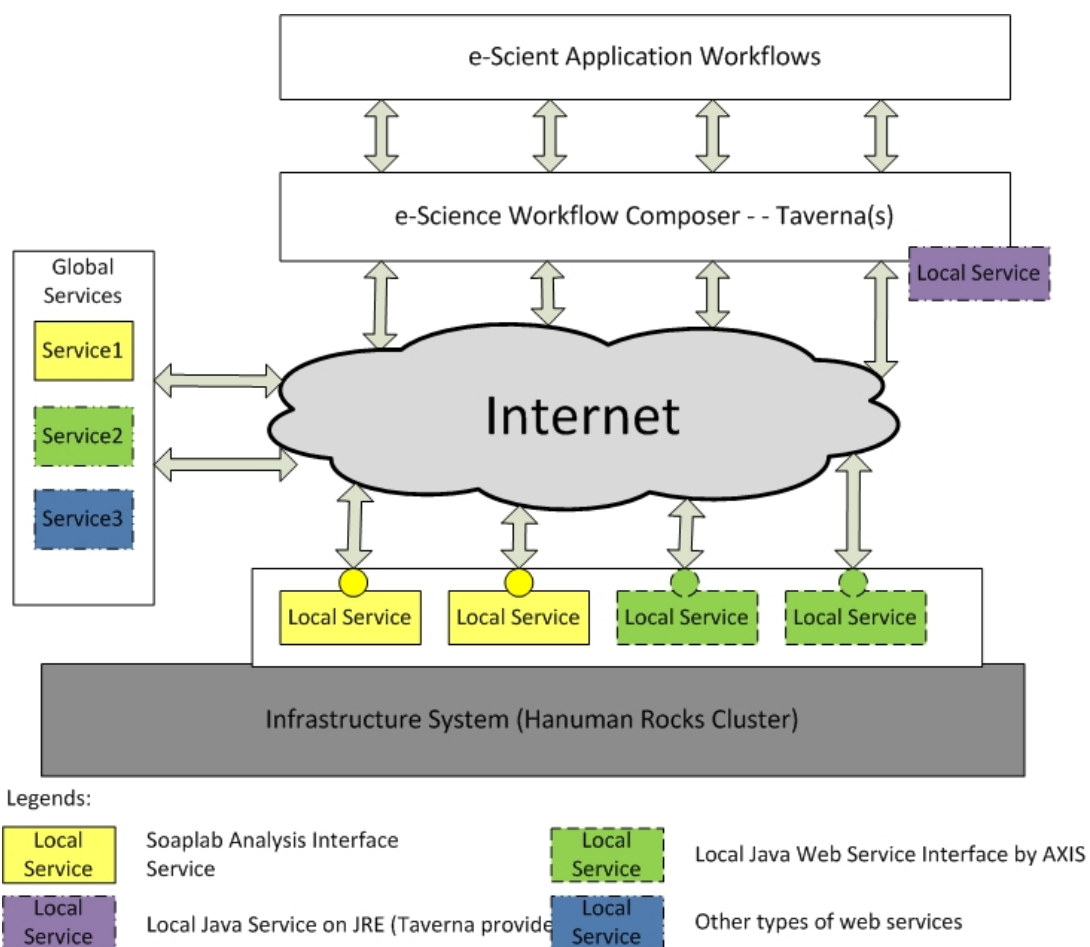
รูปที่ 3.9 ข้อความแสดงความผิดพลาดในการใช้บริการของ EBI ที่โปรแกรมทาวเวอร์นารายงาน

3.4.4 การแก้ไขปัญหาการเดินเวลาและความไม่คงเส้นคงวาในการทำงาน

วิทยานิพนธ์นี้เสนอการสร้างบริการท้องถิ่น ในส่วนการทำนายหรือการสร้างโครงสร้างต้นไม้สายวิวัฒนาการเพื่อแก้ไขปัญหาการเกินเวลา และความไม่คงเส้นคงวาในการทำงาน บริการท้องถิ่นจะสามารถช่วยประหยัดเวลาในประมวลผล และให้ผลลัพธ์ที่ถูกต้องตรงตามความต้องการมากที่สุด สถาปัตยกรรมโดยรวมของบริการท้องถิ่นที่ได้ออกแบบไว้แสดงได้ดังรูปที่ 3.10 โปรแกรมทาวเวอร์นารายงานสามารถทำงานบนเครื่องคอมพิวเตอร์ส่วนบุคคล และสามารถเข้าถึงเว็บเซอร์วิสด้วยวิธีการแบบเวิร์คโฟลว์ได้หลากหลายประเภทผ่านระบบอินเทอร์เน็ต จากรูปที่ 3.10 แสดงบริการที่เกี่ยวข้องในการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ ที่สร้างเป็นบริการท้องถิ่น ได้แก่ บริการ Soaplab แทนด้วยบล็อกสี่เหลี่ยม และบริการ Java Web Service (JWS) แทนด้วยบล็อกสี่เหลี่ยม บริการท้องถิ่นเหล่านี้จะทำงานบนระบบ Hanuman ซึ่งเป็นระบบ Rocks linux cluster

ของศูนย์กริตมหาวิทยาลัยสงขลานครินทร์ นอกจากนี้ยังจะใช้ Java local service หรือบริการท้องถิ่นที่ใช้ความสามารถของ Java interpreter ของเครื่องคอมพิวเตอร์ส่วนบุคคลของผู้ใช้งานเอง ซึ่งแสดงด้วยบล็อกสีม่วง บริการอื่นๆที่ยังสามารถทำงานได้ถูกต้อง ไม่เกิดปัญหาการเกินเวลาและความไม่คงเส้นคงวาในการทำงาน ก็ยังคงเรียกใช้งานบริการเหล่านี้จากเครือข่ายอินเทอร์เน็ต เช่นเดิม (Global services)

เราจะใช้โปรแกรม Soaplab Analysis Tool เป็นเครื่องมือสร้างตัวเชื่อมประสานโปรแกรมประยุกต์ที่ทำนายโครงสร้างต้นไม้สายวิวัฒนาการให้เป็นเว็บเซอร์วิส เพื่อให้สามารถเชื่อมต่อเข้าไปในระบบเวิร์คโฟลว์ของโปรแกรมทาเวอร์นาได้ โดยสถาปัตยกรรมของโปรแกรมทาเวอร์นารองรับการเพิ่มเว็บเซอร์วิสอยู่แล้ว



รูปที่ 3.10 สถาปัตยกรรมโดยรวมของบริการท้องถิ่น

รูปที่ 3.11 แสดงกระบวนการสร้างโครงสร้างต้นไม้สายวิวัฒนาการเพิ่มเติมในวิทยานิพนธ์นี้โดยสร้างเว็บเซอร์วิสสองชนิด คือ โปรแกรมประยุกต์สำหรับการทำนายโครงสร้างต้นไม้สายวิวัฒนาการทั้งอัลกอริทึม Parsimony สำหรับลำดับนิวคลีโอไทด์และลำดับกรดอะมิโน โดยจะห่อหุ้มหรืออิมพลีเมนต์ด้วยโปรแกรม Soaplab Analysis Tool แทนด้วยบล็อกสี่เหลี่ยม ประกอบด้วยโปรแกรม fseqboot, fdnapars, fprotpars และ fconsense ส่วนอีกชนิดคือเว็บเซอร์วิสสำหรับการจัดการอินพุตและเอาต์พุตภายในเวิร์กโฟลว์ซึ่งจะพัฒนาเป็น Java Web Service (JWS) แทนด้วยบล็อกสี่เหลี่ยม

JWS ประกอบด้วยบริการ *InputFileOperation* สำหรับทำหน้าที่ป้อนอินพุตข้อมูลลำดับนิวคลีโอไทด์เข้าสู่กระบวนการทำนายและบริการ *FecthResultOperation* สำหรับทำหน้าที่อ่านและส่งผลลัพธ์การทำงานของการทำงานทำนายโครงสร้างต้นไม้สายวิวัฒนาการ กลับมาให้โปรแกรมทาวเวอร์นาซึ่งจะได้ผลลัพธ์สองชนิดคือ โครงสร้างต้นไม้สายวิวัฒนาการที่อยู่ในรูปแบบ Treefile และอยู่ในรูปแบบ Text เช่น ตัวอย่างโครงสร้างต้นไม้สายวิวัฒนาการจากการทำนายจากการ Bootstrapping replicaton จำนวน 1,000 ชุดของลำดับนิวคลีโอไทด์ซึ่งอยู่ในรูปแบบ Treefile ดังรูปที่ 3.12 ซึ่งข้อมูลในรูปแบบ Treefile สามารถนำไปเปิดกับโปรแกรม Treeview ให้โปรแกรมสร้างเป็นโครงสร้างต้นไม้สายวิวัฒนาการ ตามประเภทที่ต้องการภายหลังได้หลายแบบได้แก่ Slanted Cladogram, Rectagular Cladogram และ Phylogram [44] ตัวอย่างโครงสร้างต้นไม้สายวิวัฒนาการจากการทำนายจากการทำ Bootstrapping replicate จำนวน 1,000 ชุดของลำดับนิวคลีโอไทด์ซึ่งอยู่ในรูปแบบ Text ดังรูปที่ 3.13

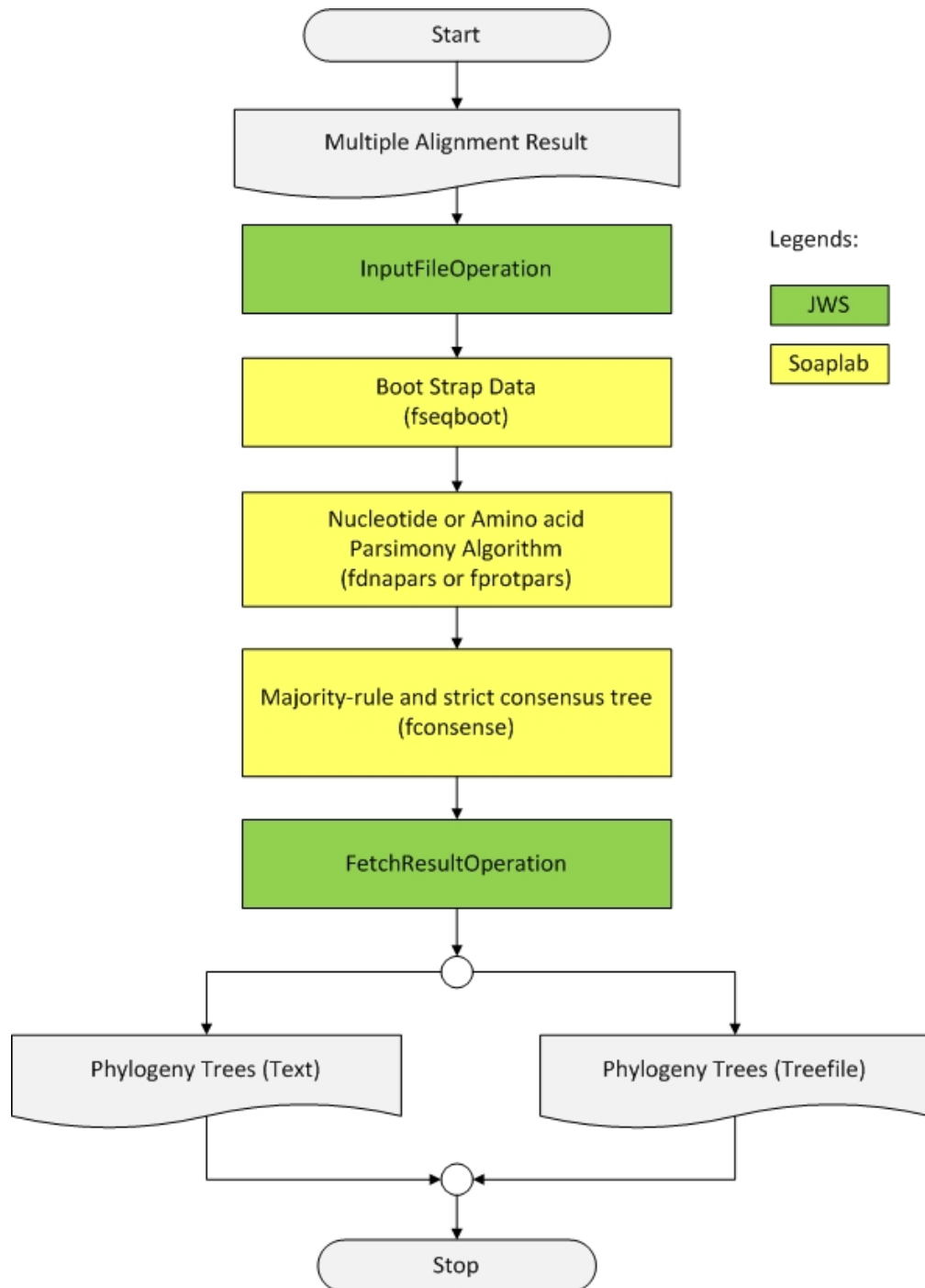
3.4.5 การพัฒนาบริการท้องถิ่น

ขั้นตอนการพัฒนาบริการท้องถิ่น สำหรับการทำนายโครงสร้างสายวิวัฒนาการของกุ่มเพื่อแก้ปัญหาการกินเวลาและความไม่คงเส้นคงวาในการทำงานงานมีดังนี้

1) การใช้งานโปรแกรม PHYLIB

โปรแกรม PHYLIB เป็นซอฟต์แวร์ชนิด Open source ของชุดโปรแกรม Emboss ซึ่งเป็นชุดเดียวกันกับที่ทำงานเป็นเว็บเซอร์วิสที่ให้บริการโดย EBI วิทยานิพนธ์นี้ติดตั้งโปรแกรม

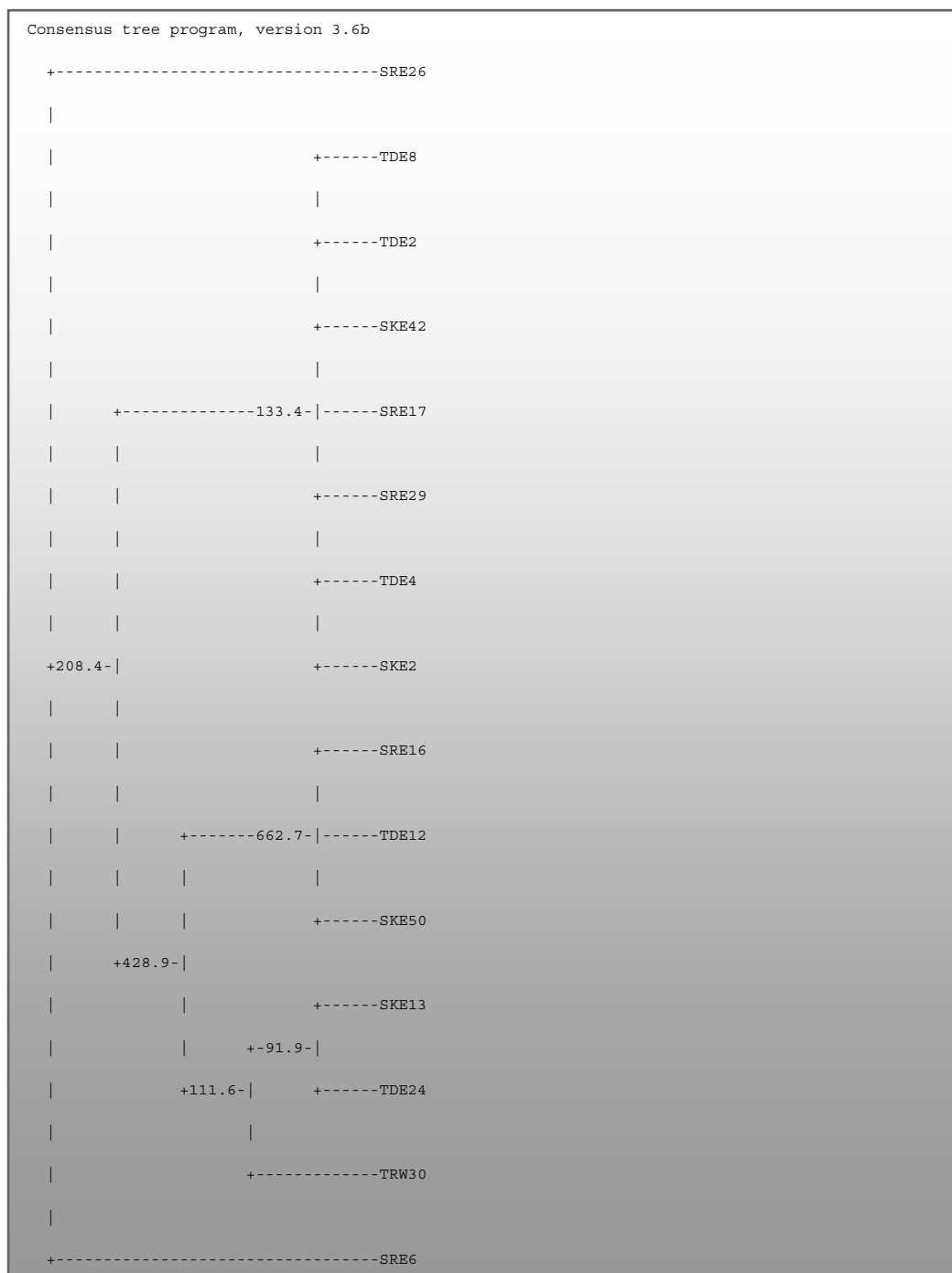
PHYLIB รุ่น 3.6b [48] สำหรับการทำนายโครงสร้างต้นไม้สายวิวัฒนาการของสิ่งมีชีวิตบนระบบคอมพิวเตอร์คลัสเตอร์ Hanuman ซึ่งใช้ระบบปฏิบัติการ Rocks Linux Cluster รุ่น 4.3 และใช้ Apache Tomcat รุ่น 5.0.28 เป็น Web service container



รูปที่ 3.11 Flowchart ของกระบวนการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ

```
(SRE26:1000.0, ((TDE8:1000.0,TDE2:1000.0,SKE42:1000.0,SRE17:1000.0,SRE29:1000.0,
TDE4:1000.0,SKE2:1000.0):133.4, ((SRE16:1000.0,TDE12:1000.0,SKE50:1000.0):662.7,
((SKE13:1000.0, TDE24:1000.0):91.9,TRW30:1000.0):111.6):428.9):208.4,SRE6:1000.0);
```

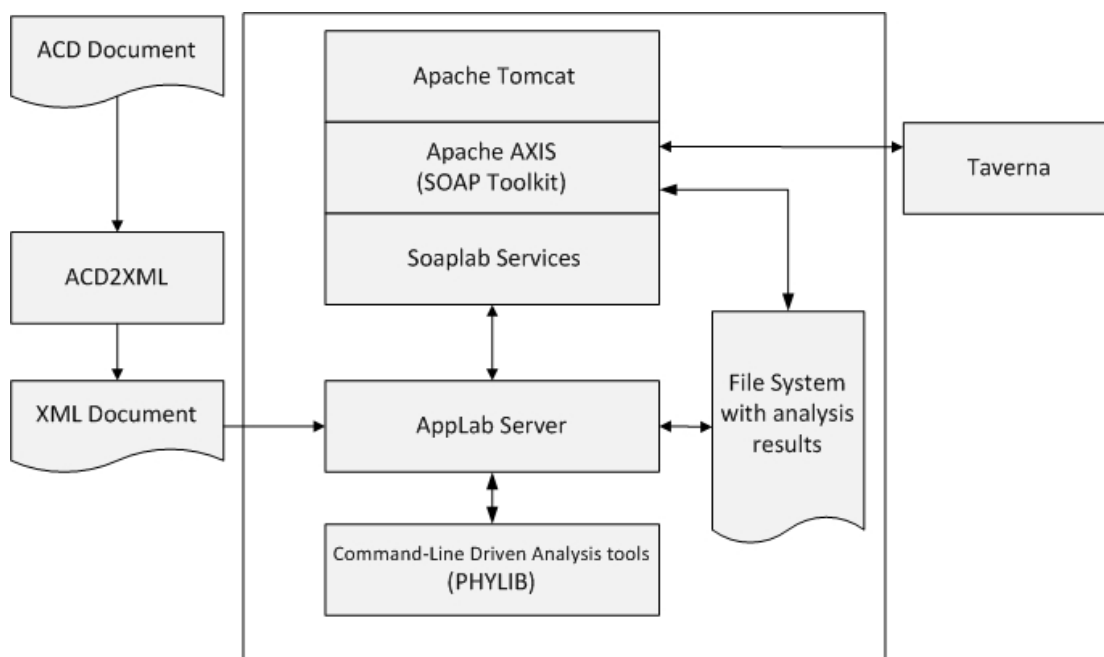
รูปที่ 3.12 โครงสร้างต้นไม้สายวิวัฒนาการในรูปแบบ Treefile ของลำดับนิวคลีโอไทด์



รูปที่ 3.13 โครงสร้างต้นไม้สายวิวัฒนาการในรูปแบบ Text ของลำดับนิวคลีโอไทด์

2) ติดตั้งโปรแกรม Soaplab Analysis Tool และ SOAP Toolkit

สถาปัตยกรรมของโปรแกรม Soaplab Analysis Tool ประกอบด้วยซอฟต์แวร์หลายชนิดประกอบกันเพื่อให้สามารถทำงานได้ตามต้องการดังที่ได้อธิบายแล้วในบทที่ 2 รูปที่ 3.14 แสดงโปรแกรม PHYLIB ในสถาปัตยกรรมการทำงานร่วมกับโปรแกรม Soaplab Analysis Tool แสดงให้เห็นถึงตำแหน่งการทำงานของโปรแกรม PHYLIB ในระบบการทำงานแบบเว็บเซอร์วิส ซึ่งจะติดต่อกับ Tool ของ Soaplab คือ AppLab Server โดย AppLab Server จะให้พอยต์เตอร์ของโปรแกรมต่างๆกับ Tool ของ Soaplab อีกโปรแกรมคือ Soaplab Services และ Tool ทั้งสองโปรแกรมทำงานบนพื้นฐานของ Apache Tomcat และ Apache AXIS ซึ่งเป็น Soap Toolkit ทำหน้าที่เชื่อมประสานการติดต่อสื่อสารออกไปสู่โลกภายนอก ในวิทยานิพนธ์นี้ก็คือระบบเวิร์คโฟลว์ในโปรแกรมทาเวอร์นา นอกจากนี้ Apache AXIS ยังทำหน้าที่เป็นโฮสต์ให้กับบริการ JWS ด้วยโปรแกรม Soaplab Analysis Tool มี Binary package ที่รวบรวม Tool เหล่านี้ไว้เรียบร้อยแล้วคือ Analysis Interface [37]



รูปที่ 3.14 สถาปัตยกรรมการทำงานของโปรแกรมต่างๆของบริการท้องถิ่น

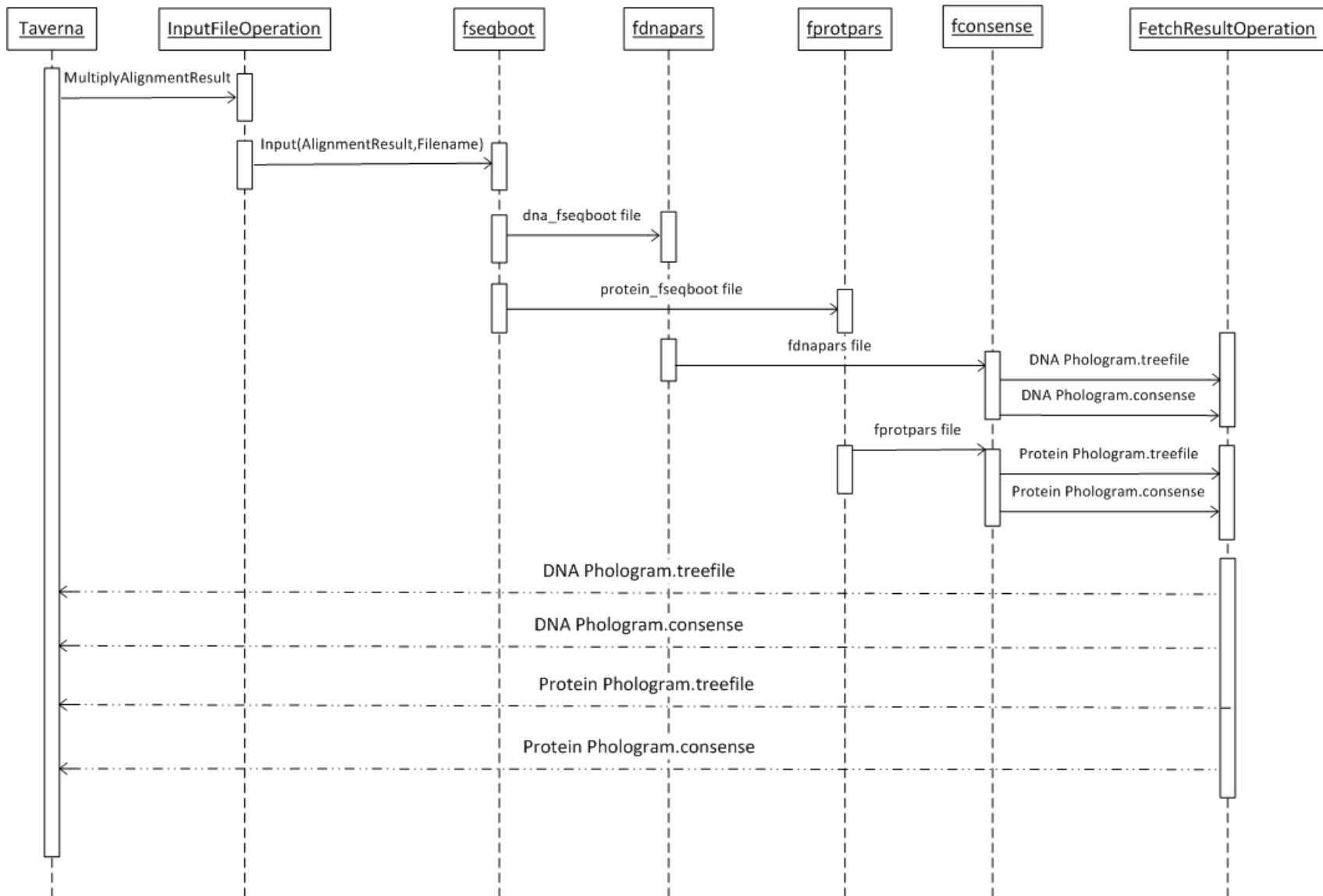
3) การสร้างตัวเชื่อมประสานเว็บเซอร์วิส (Web service interface)

การสร้างตัวเชื่อมประสานเว็บเซอร์วิสหรือ Web service interface ให้กับชุดโปรแกรมในการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ สามารถสร้างได้ตามสถาปัตยกรรมการทำงานในรูปแบบที่ 3.14 เครื่องมือของ Soaplab Analysis Tool ที่สำคัญอีกตัวหนึ่งคือ ACD2XML ซึ่งทำหน้าที่ในการแปล (translate) เอกสาร AJAX Command Definition (ACD) [49] เป็นเอกสาร XML เนื่องจากเอกสาร XML เป็นภาษามาตรฐานของโปรโตคอล SOAP ซึ่งเป็นหัวใจของการทำงาน ดังนั้นขั้นตอนสำคัญคือ การสร้างเอกสาร ACD เพื่ออธิบายการทำงานและการปรับแต่งของชุดโปรแกรม PHYLIB ที่อยู่ในรูปแบบ Interactive command line ให้เข้ากับการใช้งานในงานวิทยานิพนธ์คือ Non-interactive command line จากนั้นจึงใช้ ACD2XML แปลเอกสาร ACD เป็นเอกสาร XML เพื่อส่งต่อเข้าสู่กระบวนการทำงานของ Soaplab Analysis Tool ต่อไป

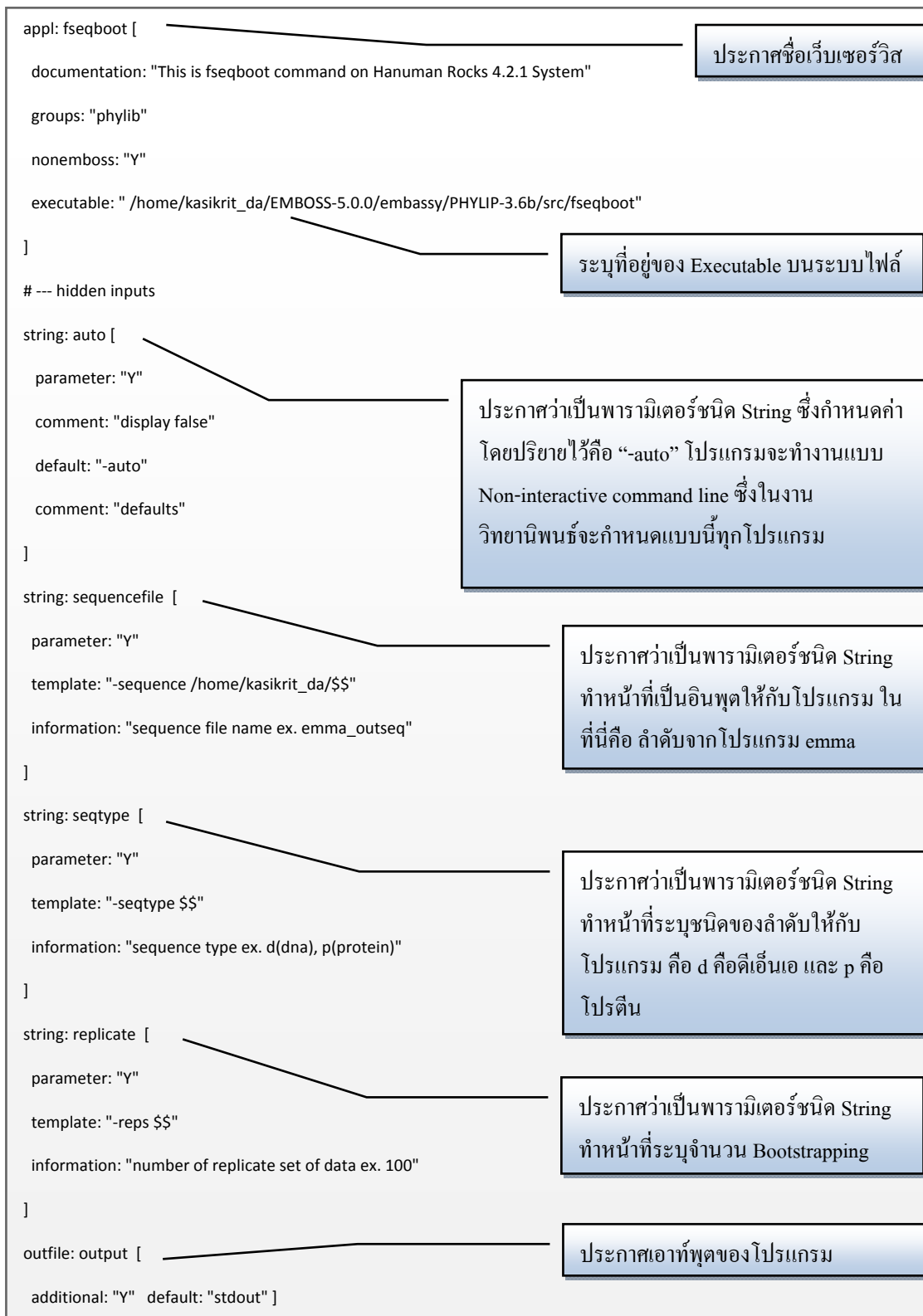
ชุดโปรแกรม PHYLIB สำหรับการทำงานโครงสร้างสายต้นไม้วิวัฒนาการ ประกอบด้วยโปรแกรม fseqboot, fdnapars, fprotpars และ fconsense ซึ่งลำดับการทำงานภายในเวิร์กโฟลว์ของโปรแกรมทาวเวอร์นาและการส่งต่อข้อมูลสามารถอธิบายด้วย Sequence Diagram ดังรูปที่ 3.15 โดยประกอบด้วยการทำงานด้วยเอกสาร ACD ดังต่อไปนี้

- เอกสาร ACD สำหรับโปรแกรม fseqboot

โปรแกรม fseqboot สามารถรับอินพุตได้ทั้งลำดับนิวคลีโอไทด์และลำดับกรดอะมิโน ซึ่งการแยกประเภทของลำดับกำหนดไว้เป็นพารามิเตอร์ ยิ่งกำหนดให้โปรแกรม fseqboot ทำ Bootstrapping replicate มากเท่าไร ก็จะได้โครงสร้างต้นไม้สายวิวัฒนาการที่มีความเป็นไปได้ และถูกต้องมากยิ่งขึ้น ผลลัพธ์ที่ได้จะส่งต่อให้โปรแกรม fdnapars สำหรับวิเคราะห์ลำดับนิวคลีโอไทด์ และ fprotpars สำหรับวิเคราะห์ลำดับกรดอะมิโน โดยในวิทยานิพนธ์นี้จะทดสอบจำนวน Bootstrapping replicate สูงสุดถึง 50,000 ชุดเพื่อการทดสอบประสิทธิภาพของบริการท้องถิ่นเอกสาร ACD สำหรับโปรแกรม fseqboot ดังรูปที่ 3.16



รูปที่ 3.15 Sequence diagram ของการทำงานภายในเวิร์คโฟลว์ของโปรแกรมทาวเอร์นาในการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ



รูปที่ 3.16 เอกสาร ACD สำหรับโปรแกรม fseqboot

- เอกสาร ACD สำหรับโปรแกรม fdnapars

โปรแกรม fdnapars มีอัลกอริทึมการทำงาน DNA Parsimony สำหรับการเตรียมข้อมูลเพื่อส่งต่อให้กับโปรแกรม fconsense ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์ เอกสาร ACD สำหรับโปรแกรม fdnpars ดังรูปที่ 3.17

```

appl: fdnapars [
  documentation: "This is fdnapars command on Hanuman Rocks 4.2.1 System"
  groups: "phylip"
  nonemboss: "Y"
  executable: "/home/kasikrit_da/EMBOSS-5.0.0/embassy/PHYLIP-3.6b/src/fdnapars"
]
# --- hidden inputs
string: auto [
  parameter: "Y"
  comment: "display false"
  default: "-auto"
  comment: "defaults"
]
string: fseqbootfile [
  parameter: "Y"
  template: "-sequence /home/kasikrit_da/analysis-interfaces/a/unknown/Projects/default/Data/$$fseqboot"
  information: "sequence file name ex. emma_outseq"
]
outfile: output [
  additional: "Y"
  default: "stdout"
]

```

ระบุที่อยู่ของ Executable บนระบบไฟล์

ระบุที่อยู่ของอินพุตให้กับ โปรแกรม fdnpars ซึ่งก็คือเอาที่พุดของ โปรแกรม fseqboot ของ ลำดับนิวคลีโอไทด์

รูปที่ 3.17 เอกสาร ACD สำหรับโปรแกรม fdnpars

- เอกสาร ACD สำหรับโปรแกรม fprotpars

โปรแกรม fprotpars มีอัลกอริทึมการทำงาน Protein parsimony สำหรับการเตรียมข้อมูลเพื่อส่งต่อไปให้กับโปรแกรม fconsense ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโนต่อไป เอกสาร ACD สำหรับโปรแกรม fprotpars ดังรูปที่ 3.18

```

appl: fprotpars [
  documentation: "This is fprotpars command on Hanuman Rocks 4.2.1 System"
  groups: "phylip"
  nonemboss: "Y"
  executable: "/home/kasikrit_da/EMBOSS-5.0.0/embassy/PHYLIIP-3.6b/src/fprotpars"
]
# --- hidden inputs
string: auto [
  parameter: "Y"
  comment: "display false"
  default: "-auto"
  comment: "defaults"
]
string: fseqbootfile [
  parameter: "Y"
  template: "-sequence /home/kasikrit_da/analysis-interfaces/a/unknown/Projects/default/Data/$$fseqboot"
  information: "sequence file name ex. emma_outseq"
]
outfile: output [
  additional: "Y"
  default: "stdout"
]

```

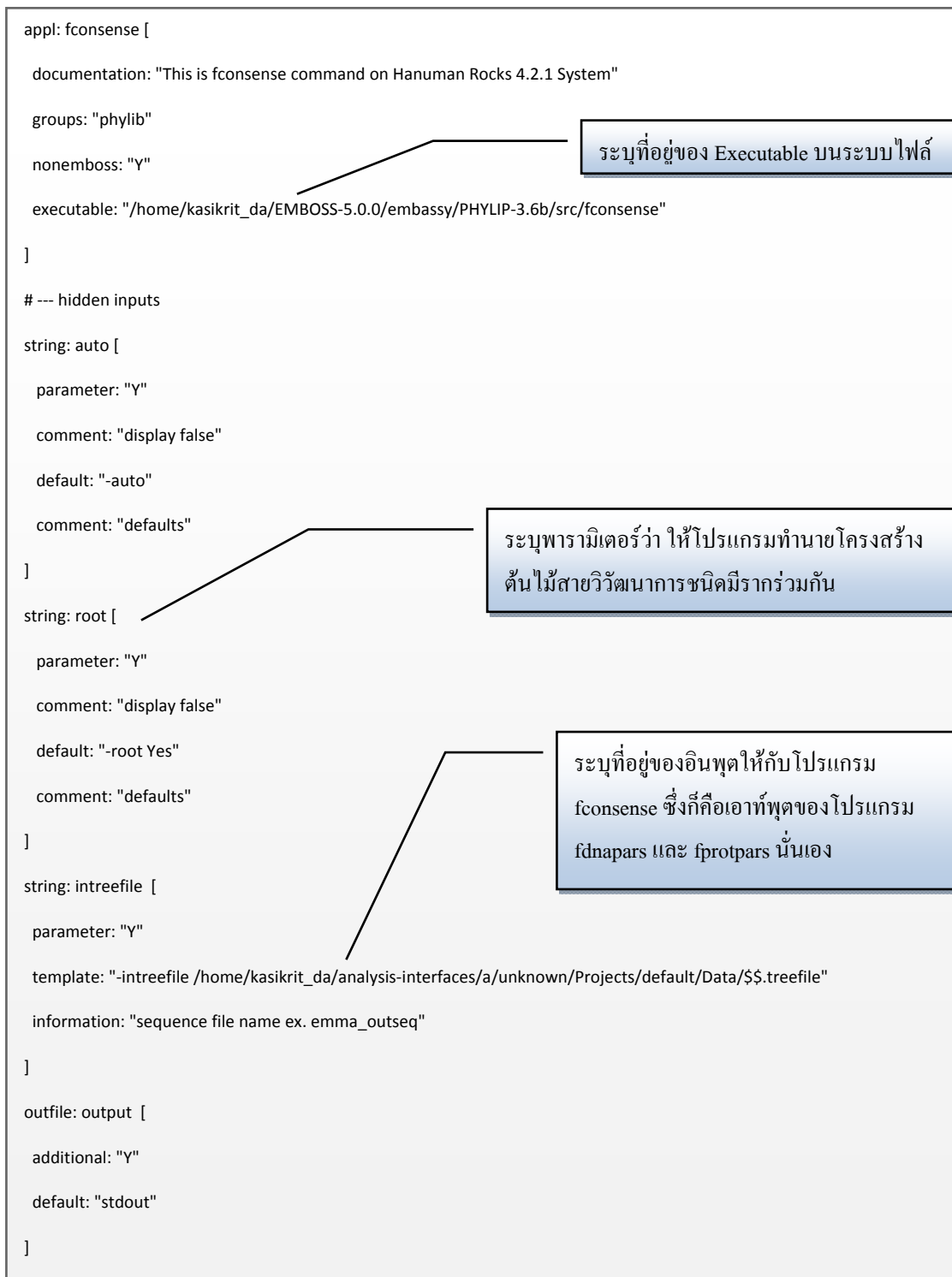
ระบุที่อยู่ของ Executable บนระบบไฟล์

ระบุที่อยู่ของอินพุตให้กับโปรแกรม fdnpars ซึ่งก็คือเอาท์พุตของโปรแกรม fseqboot ของลำดับกรดอะมิโน

รูปที่ 3.18 เอกสาร ACD สำหรับโปรแกรม fprotpars

- เอกสาร ACD สำหรับโปรแกรม fconsense

โปรแกรม fonsense สามารถรับอินพุตได้ทั้งลำดับนิวคลีโอไทด์และลำดับกรดอะมิโนซึ่งสามารถแยกแยะได้โดยการระบุพารามิเตอร์ โปรแกรม fconsense ใช้ทำนายโครงสร้างต้นไม้สายวิวัฒนาการที่มีความเป็นไปได้มากที่สุด เอกสาร ACD สำหรับโปรแกรม fconsense ดังรูปที่ 3.19 และให้อาท์พุตเป็นข้อมูลโครงสร้างต้นไม้สายวิวัฒนาการ 2 รูปแบบคือ Treefile และ Text ดังที่ได้กล่าวแล้ว



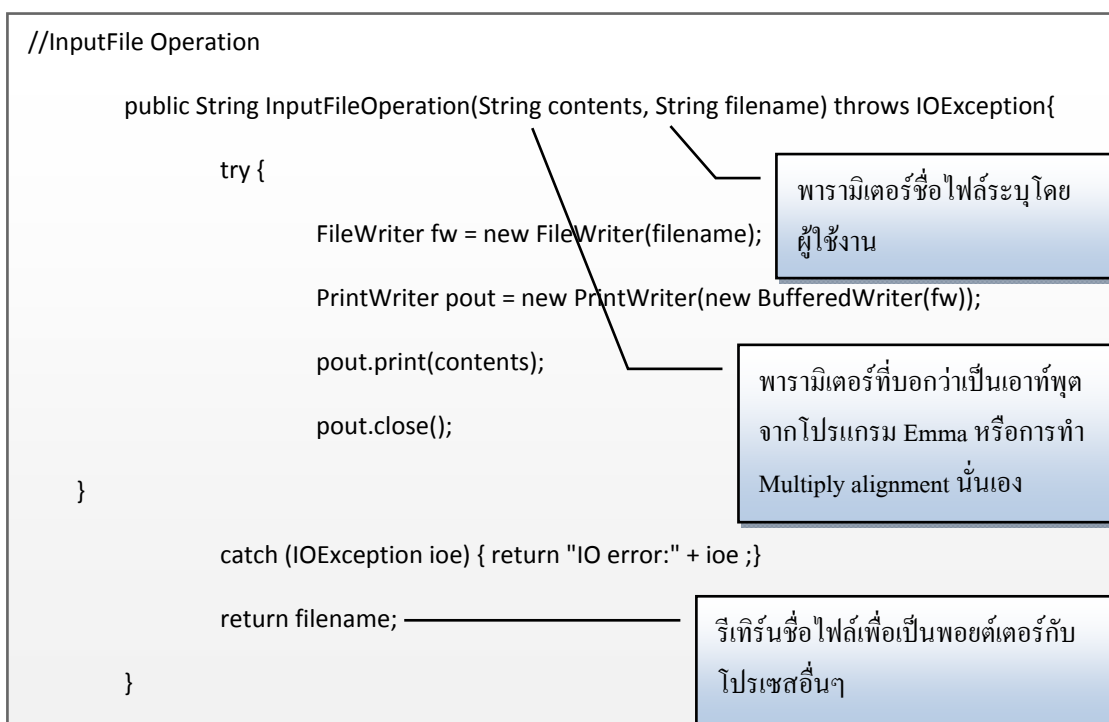
รูปที่ 3.19 เอกสาร ACD สำหรับโปรแกรม fconsense

4) สร้างบริการท้องถิ่นสำหรับการจัดการด้านอินพุตและเอาต์พุต

บริการท้องถิ่นที่สร้างในงานวิทยานิพนธ์เป็นชนิด JWS ซึ่งเป็น Java class ทำหน้าที่จัดการด้านอินพุตและเอาต์พุตของการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ บนระบบไฟล์ ตามที่ได้กำหนดไว้ในเอกสาร ACD ส่วนของ Java Class ที่เป็นแกนสำคัญของการทำงานเป็นดังนี้

- บริการท้องถิ่นสำหรับการจัดการอินพุต

บริการท้องถิ่นสำหรับการจัดการอินพุต ทำหน้าที่ป้อนข้อมูลอินพุตซึ่งเป็นลำดับนิวคลีโอไทด์และโปรตีนที่เป็นผลลัพธ์จากโปรแกรม Emma ของสถาบัน EBI ซึ่งทำงานอยู่บนเครือข่ายอินเทอร์เน็ตเข้าสู่กระบวนการทำนายโครงสร้างต้นไม้สายวิวัฒนาการ ซึ่งทำงานในระบบเวิร์กโฟลว์ของโปรแกรมทาวเวอร์นาดังรูปที่ 3.20



รูปที่ 3.20 Java class ของบริการท้องถิ่นสำหรับการจัดการอินพุต

- บริการท้องถิ่นสำหรับการจัดการเอาท์พุต

บริการท้องถิ่นสำหรับการจัดการเอาท์พุต ทำหน้าที่อ่านผลลัพธ์ของการทำนาย โครงสร้างสายวิวัฒนาการทั้งของลำดับนิวคลีโอไทด์และโปรตีน จากระบบไฟล์แล้วส่งผลลัพธ์ กลับมาให้ที่ระบบเวิร์กโฟลว์ในโปรแกรมทาวเวอร์นาคังรูปที่ 3.21 และ รูปที่ 3.22

1) Java class สำหรับการอ่านและส่งผลลัพธ์ของรูปแบบ Treefile

```
//Reading a file and show content of emboss result in analysis_interface dir

public String ReadFileEmbossTreefile(String filename)throws IOException{

    String s, content="";

    try {

        FileReader fin = new FileReader("/home/kasikrit_da/analysis-
        interfaces/a/unknown/Projects/default/Data/"+filename+".treefile");

        BufferedReader bin = new BufferedReader(fin);

        while( (s = bin.readLine()) != null){

            content += s + "\n";

        }

        bin.close();

    }

    catch (IOException ioe) {

        System.out.println( "IO error:" + ioe);

    }

    return content;

}

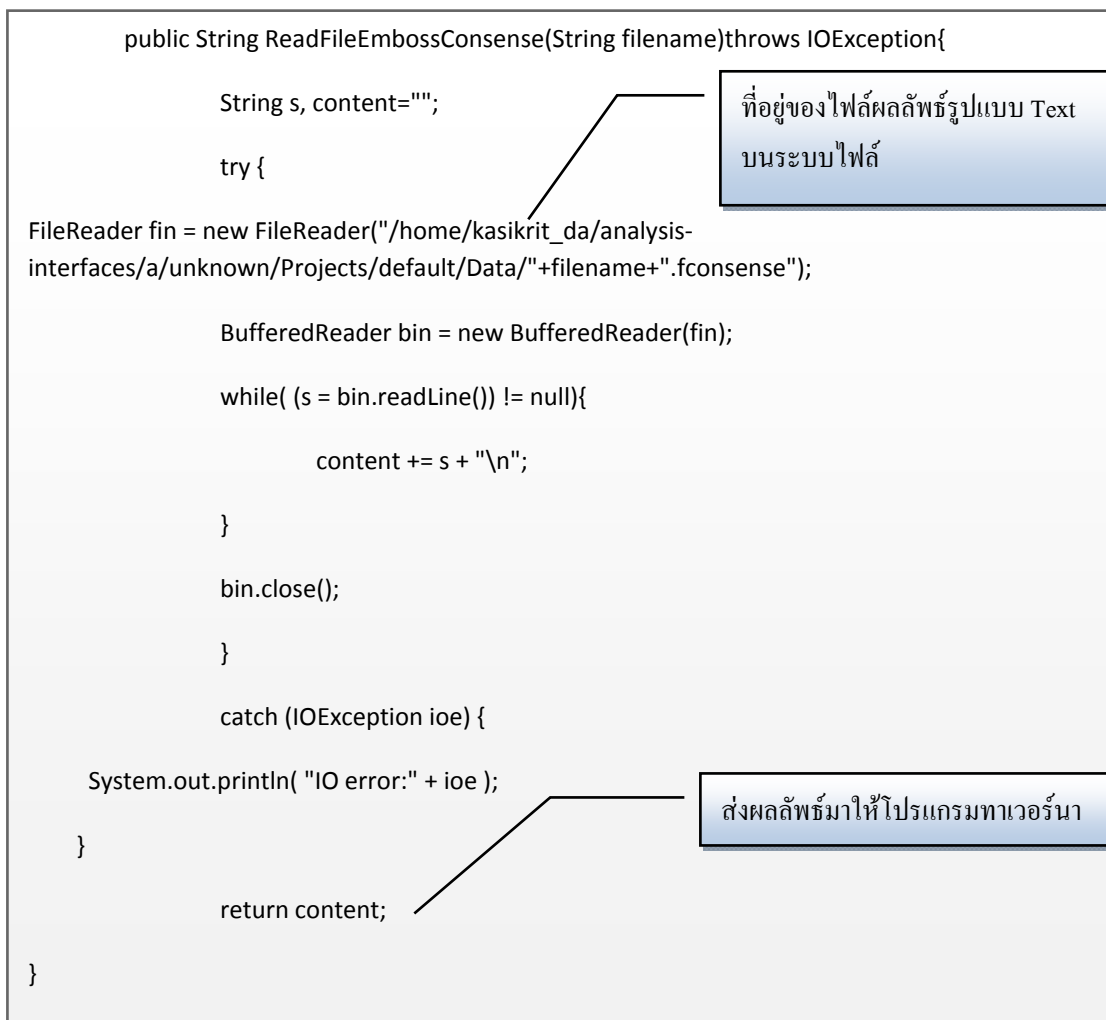
}
```

ที่อยู่ของไฟล์ผลลัพธ์รูปแบบ Treefile บนระบบไฟล์

ส่งผลลัพธ์มาให้โปรแกรมทาวเวอร์นา

รูปที่ 3.21 Java class สำหรับการอ่านและส่งผลลัพธ์ของรูปแบบ Treefile

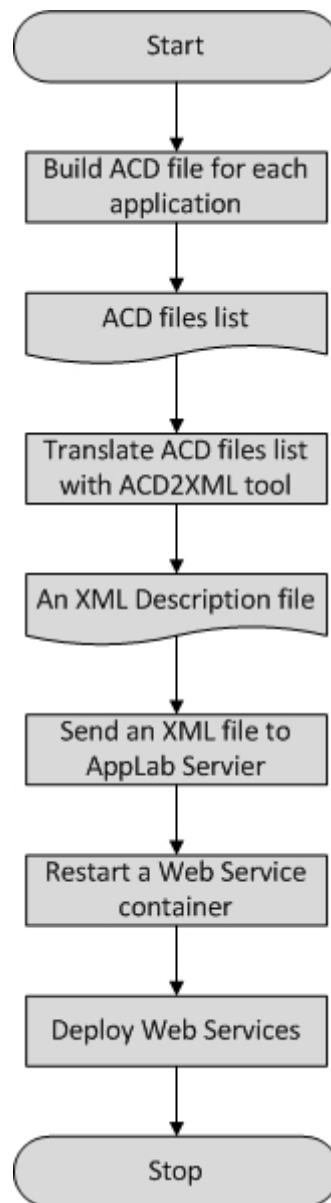
2) Java class สำหรับการอ่านและส่งผลลัพธ์ในรูปแบบ Text



รูปที่ 3.22 Java class สำหรับการอ่านและส่งผลลัพธ์ในรูปแบบ Text

5) การใช้งานบริการท้องถิ่น

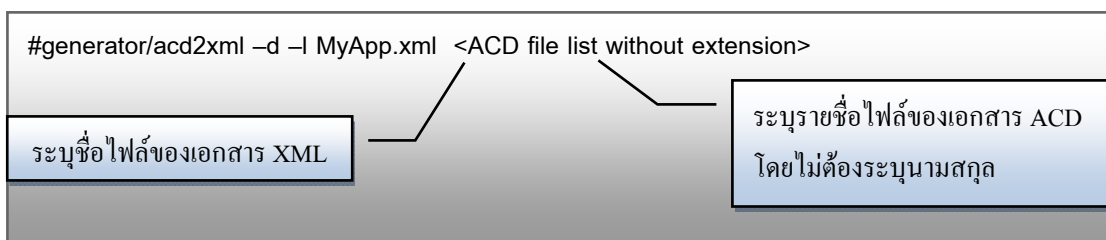
การใช้งานบริการท้องถิ่น (Local web services employing) ด้วยโปรแกรม Soaplab Analysis Tool โดยมี Flowchart แสดงขั้นตอนดังรูปที่ 3.23 และอธิบายการใช้งานบริการท้องถิ่นดังนี้



รูปที่ 3.23 ขั้นตอนการใช้งานบริการท้องถิ่นด้วยโปรแกรม Soaplab Analysis Tool

5.1) เอกสาร ACD จะต้องประกอบด้วยนามสกุล acd ได้แก่ fseqboot.acd, fdnapars.acd, fprotpars.acd และ fconsense.acd จากนั้นบันทึกเอกสาร ACD ที่ระบบไฟล์ soaplab/src/etc/acd/applab หรือ analysis-interfaces/metadata สำหรับชุดโปรแกรมแบบ Binary package

5.2) แปลงเอกสาร ACD เป็นเอกสาร XML ด้วยเครื่องมือ ACD2XML ดังรูปที่



รูปที่ 3.24 การใช้เครื่องมือ ACD2XML

5.3) ส่งเอกสาร XML ให้กับ AppLab Server โดยให้ทำงานเป็นโปรเซสเบื้องหลัง

```
#run-AppLab-server MyApp.xml &
```

5.5) รีสตาร์ทโปรเซสของ Apache Tomcat

5.6) เรียกใช้งานเซอร์วิสโดยให้ทำงานเป็นโปรเซสเบื้องหลัง

```
#!/ws/deploy-web-services -d -A MyApp.xml &
```

5.7) เอ็นพอยต์ของเว็บเซอร์วิสคือ `http://<hostname>:<port>/axis/services` ในวิทยานิพนธ์คือ `http://hanuman.psu.ac.th:8081/axis/services` ซึ่งสามารถตรวจสอบการเข้าถึงกับเว็บเบราว์เซอร์และโปรแกรมทาวเวอร์นา

5.8) สำหรับบริการท้องถิ่นชนิด JWS สามารถเรียกใช้งานเซอร์วิสได้โดยบันทึก Java Class เป็นนามสกุล .jws แล้วไว้ที่ระบบไฟล์ `<tomcat dir>/webapps/axis` ในวิทยานิพนธ์นี้คือ `tomcat-5.0.28/webapps/axis` และเอ็นพอยต์ในการเข้าถึงบริการคือ `http://<hostname>:<port>/axis/<java class name>.jws?wsdl` ในที่นี้คือ `http://hanuman.psu.ac.th:8081/axis /HanumanWSDL.jws?wsdl`

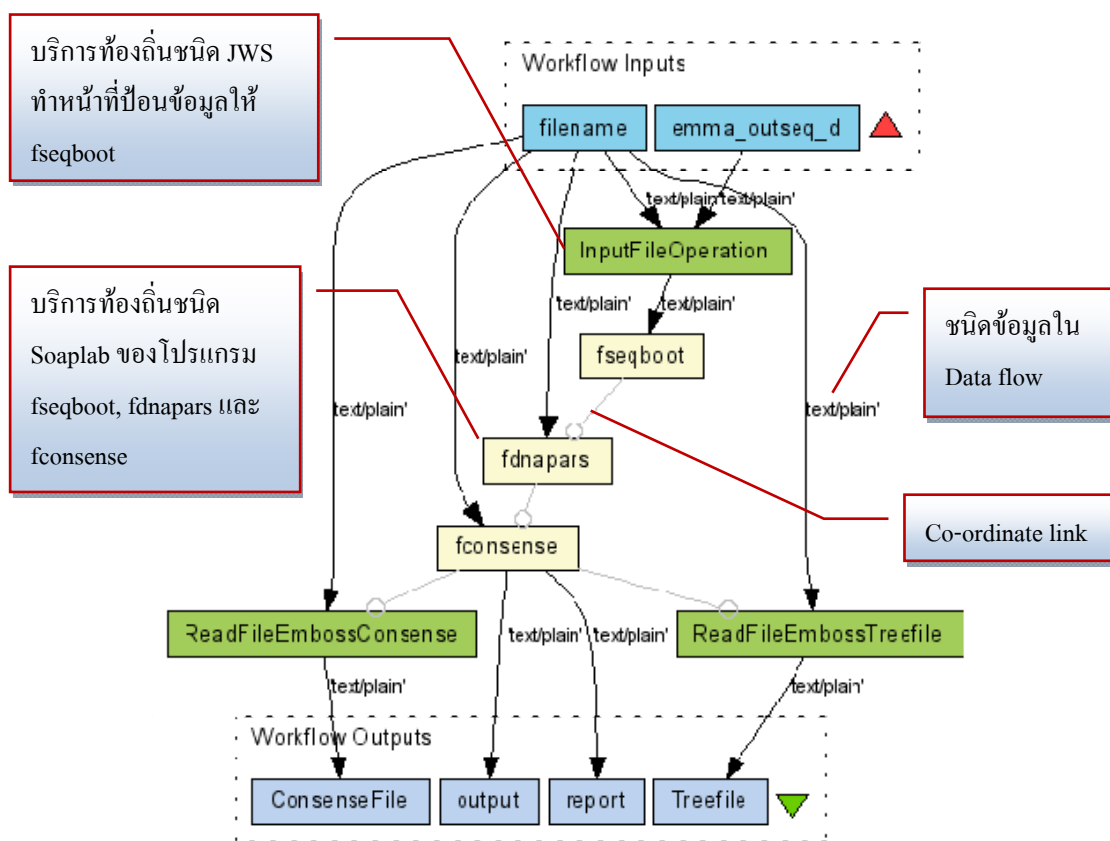
5.9) การรีสตาร์ทโปรเซสของ Apache Tomcat หลังจากกระบวนการการเรียกใช้งานเซอร์วิสสามารถทำได้โดยอิสระ ไม่กระทบการทำงานของบริการต่อที่ทำงานเป็นโปรเซสเบื้องหลัง แต่การเรียกใช้งานเซอร์วิสใหม่ จำเป็นต้องรีสตาร์ทโปรเซสของ Apache Tomcat ทุกครั้งเพื่อให้ระบบรับรู้ถึงบริการใหม่ที่เพิ่มเข้ามา

3.4.6 การพัฒนาเวิร์คโฟลว์ที่ใช้บริการท้องถิ่น

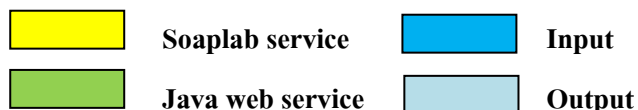
บริการท้องถิ่นที่สร้างขึ้นทั้งบริการ Soaplab และบริการ JWS สามารถเข้าถึงได้ด้วยโปรแกรมทาวเวอร์นาและสามารถสร้าง Instance ของบริการเข้าไปในระบบเวิร์คโฟลว์เพื่อการทำนายโครงสร้างต้นไม้สายวิวัฒนาการได้ ดังเวิร์คโฟลว์ที่จะกล่าวถึงต่อไปนี้

- การพัฒนาเวิร์กโฟลว์ของแต่ละกระบวนการทำนาย

รูปที่ 3.25 แสดงเวิร์กโฟลว์ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์เพราะมีการใช้งานบริการท้องถิ่น fdnapars ในการทำงานอัลกอริทึม DNA parsimony สังเกตได้ว่ามีการใช้ Co-ordinate link หรือการเชื่อมโยงแบบควบคุม คือบริการถัดไปจะเริ่มทำงานได้ก็ต่อเมื่อบริการก่อนหน้าทำงานเสร็จสิ้นเสียก่อน

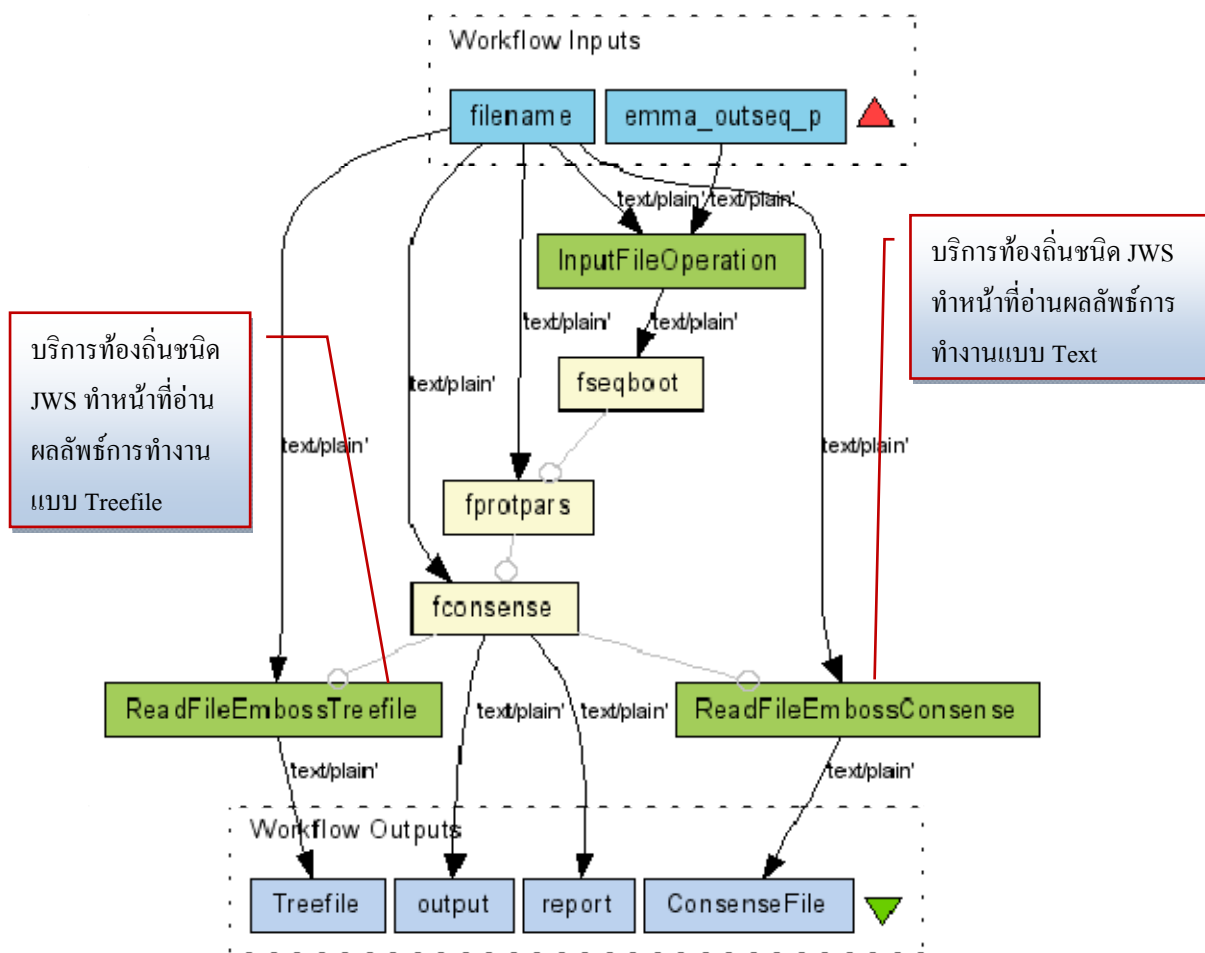


Legends:



รูปที่ 3.25 เวิร์กโฟลว์ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์

รูปที่ 3.26 เวิร์กโฟลว์ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโน เพราะมีการใช้งานบริการท้องถิ่น fprotpars ในการทำงานอัลกอริทึม Protein parsimony จะสังเกตว่ามีการใช้ Co-ordinate link หรือการเชื่อมโยงแบบควบคุมเช่นเดียวกัน



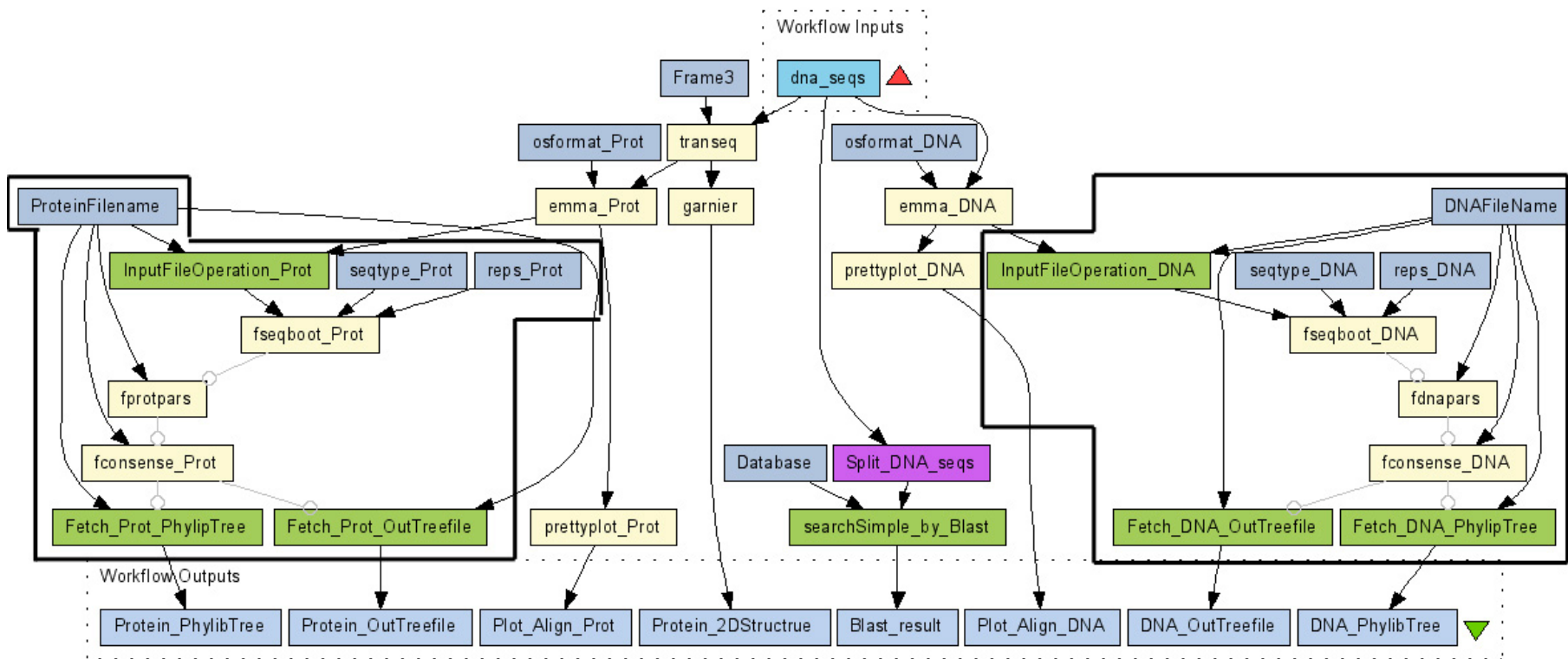
Legends:

- Soaplab service
- Input
- Java web service
- Output

รูปที่ 3.26 เวิร์กโฟลว์ทำนายโครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโน

- การพัฒนาเวิร์คโฟลว์โดยรวมระบบงานทั้งหมดเข้าด้วยกัน

จากการทดสอบการทำงานอย่างง่ายก็ทำให้ทราบว่าเวิร์คโฟลว์ที่สร้างในรูปที่ 3.25 และ รูปที่ 3.26 มีลำดับกระบวนการทำงาน การเข้าจังหวะการทำงานและมี Data flow ไหลในเวิร์คโฟลว์ได้ถูกต้องสมบูรณ์ดีแล้ว จึงทำการปรับแก้เวิร์คโฟลว์หลักเวิร์คโฟลว์แรก (รูปที่ 3.7) โดยเปลี่ยนมาใช้บริการท้องถิ่นที่ได้ทำการสร้างขึ้นมาใช้งานเอง จึงได้เวิร์คโฟลว์สำหรับการทำนายโครงสร้างต้นไม้สายวิวัฒนาการของกิ้งก่าใช้บริการท้องถิ่น (ในกรอบสีดำ) ดังรูปที่ 3.27



Legends:

- | | | | |
|---|---|--|--|
|  Input |  EBI Soaplab Web Service |  Standard WSDL Service |  Co-ordinate |
|  Parameter for service |  Shim Service |  Output |  Dataflow between service |

รูปที่ 3.27 เวิร์กโฟลว์ใหม่ที่ใช้บริการท้องถิ่น (ในกรอบสีดำ)

3.4.7 ผลการทดลอง

ตารางที่ 3.5 แสดงผลการทดสอบของกรณีทดสอบต่างๆของเวิร์คโฟลว์ใหม่ที่ใช้บริการท้องถิ่นในรูปแบบที่ 3.27 โดยใช้จำนวน Bootstrapping replicate ของลำดับนิวคลีโอไทด์และกรดอะมิโนโดยจะเพิ่มขึ้นเรื่อยๆจนกระทั่งถึง 50,000 ชุด เวลาที่เวิร์คโฟลว์ทำงานแสดงในหน่วยนาฬิกาซึ่งได้จากรายงานสถานะการทำงานจากโปรแกรมทาเวอร์น่า

ตารางที่ 3.5 กรณีทดสอบและผลการทดสอบของเวิร์คโฟลว์ใหม่ที่ใช้บริการท้องถิ่น

จำนวนชุดข้อมูลของการทำซ้ำ (No. of Replicates)	เวลาที่ใช้ในการทำงาน (นาฬิกา) (Time Used)
100	02.12
250	02.06
500	02.06
1,000	02.06
2,000	02.24
3,000	03.04
4,000	03.47
5,000	03.52
10,000	08.09
20,000	15.08
30,000	22.02
40,000	29.22
50,000	36.17

ตารางที่ 3.6 แสดงการเปรียบเทียบกรณีทดสอบและผลการทดสอบของวิธีการทำงานแบบต่างๆของลำดับนิวคลีโอไทด์คือ วิธีการแบบคัดลอกแล้ววาง (Manual) กับวิธีอัตโนมัติ

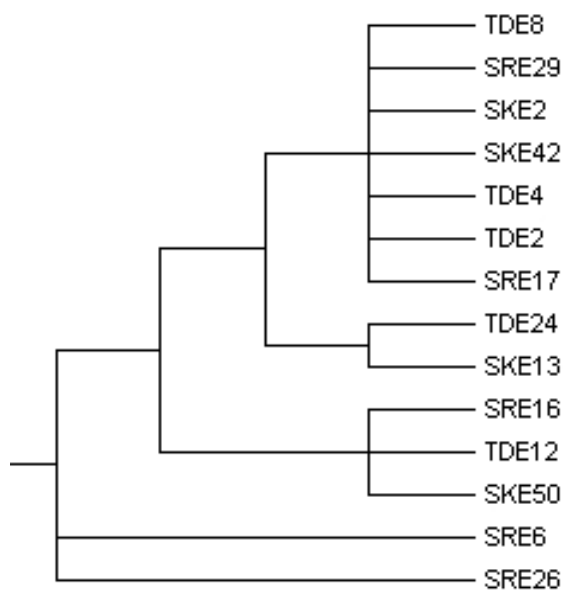
ด้วยการใช้เวิร์คโฟลว์ โดยมีสองประเภทคือ วิธีการแบบเวิร์คโฟลว์ที่ไม่ได้ใช้บริการท้องถิ่น และวิธีการแบบเวิร์คโฟลว์โดยใช้บริการท้องถิ่น สำหรับการทำให้ Bootstrapping replicate 1,000 ชุด เวิร์คโฟลว์แรกใช้เวลาทำงาน 2.06 นาที แต่ล้มเหลวในการทำงานเนื่องจากสาเหตุการเกินเวลาดังที่กล่าวแล้วที่จำนวนชุดข้อมูลของการทำให้ Bootstrapping replicate 2,000 ชุด วิธีการแบบคัดลอกและวางโดยนักชีวสารสนเทศใช้เวลาประมาณ 15 นาที ในขณะที่วิธีการทำงานแบบอัตโนมัติด้วยเวิร์คโฟลว์ที่ให้บริการท้องถิ่นใช้เวลา 2.24 นาที ดังนั้นการลดลงของเวลาในการประมวลผลไปเป็นอย่างมากมีนัยสำคัญ เมื่อเราเพิ่มตัวเลขการทำ Bootstrapping replicate ของลำดับนิวคลีโอไทด์เป็น 20,000 และ 50,000 ชุด เวิร์คโฟลว์ที่ให้บริการท้องถิ่นใช้เวลาทำงาน 15.08 และ 36.17 นาทีตามลำดับ วิทยานิพนธ์นี้ไม่ได้ทดสอบกรณีทดสอบนี้โดยวิธีการคัดลอกแล้ววางเนื่องจากทดสอบทำได้ยากมากหรือเป็นไปได้ที่จะคัดลอกแล้ววางข้อมูลจำนวนมากเหล่านี้ ระหว่างเว็บเพจของกระบวนการต่างๆ

ตารางที่ 3.6 ตารางการเปรียบเทียบกรณีทดสอบและผลการทดสอบของวิธีการทำงานแบบต่างๆ

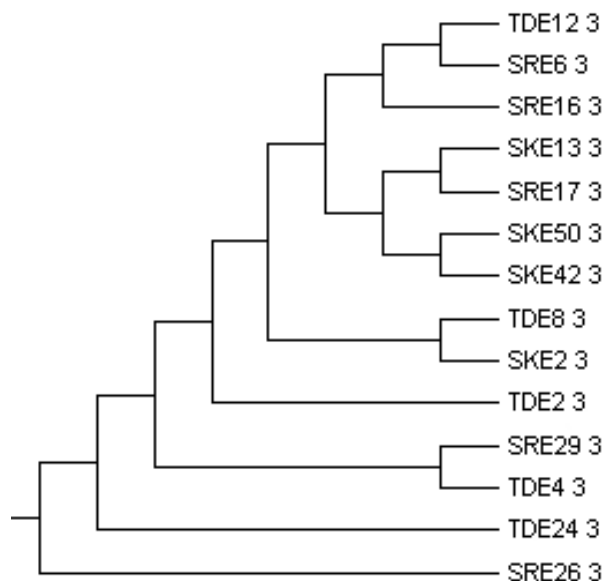
จำนวนชุดข้อมูลของการทำซ้ำ (No. of Replicates)	กรณีทดสอบ (Test Case)			Running Time (Min.)
	Manual	Automatic Workflow		
		No local service	With local web services	
1,000	-	✓	-	02.06
2,000	✓	-	-	15
2,000	-	✓	-	ล้มเหลว
2,000	-	-	✓	2.24
20,000	✓	-	-	ยากหรือเป็นไปได้ที่จะทดสอบ
20,000	-	-	✓	15.08
50,000	-	-	✓	30.17

รูปที่ 3.28 และรูปที่ 3.29 คือโครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์และกรดอะมิโนตามลำดับ ซึ่งได้รับการตรวจสอบความถูกต้องจากศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์ เป็นความสัมพันธ์ทางระบบวิวัฒนาการของกิ่งจากแหล่งที่อยู่อาศัยต่างๆกัน 3 แหล่ง ความรู้ที่ได้จะสามารถช่วยคัดเลือกสายพันธุ์ในการศึกษาด้านอื่นๆต่อไปในอนาคต

โครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์ในรูปที่ 3.28 แสดงความสัมพันธ์ทางด้านพันธุกรรม อย่างเช่น ลำดับTDE24 และลำดับ SKE13 อยู่ในกลุ่มประเภท (Species) ด้วยกันและมีความสัมพันธ์ใกล้ชิดกับลำดับ SRE17 มากกว่าลำดับ SRE16 เป็นต้น โครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโนในรูปที่ 3.29 แสดงให้เห็นว่าลำดับ SKE50 และลำดับ SKE42 อยู่ในกลุ่มประเภทเดียวกันและต่างก็สืบสายพันธุ์มาจาก TDE8 เป็นต้น



รูปที่ 3.28 โครงสร้างต้นไม้สายวิวัฒนาการของลำดับนิวคลีโอไทด์ของกุ้ง



รูปที่ 3.29 โครงสร้างต้นไม้สายวิวัฒนาการของลำดับกรดอะมิโนของกุ้ง

3.5 อภิปรายผลการทดลอง

วิธีการทำงานแบบเก่าหรือคัดลอกแล้ววางสามารถทำงานกับจำนวน 2,000 Bootstrapping replicates ได้ โดยนักชีวสารสนเทศที่มีความชำนาญในระบบงานและการใช้โปรแกรมแบบออนไลน์กับเว็บแอปพลิเคชัน ใช้เวลาในการทำงานประมาณ 15 นาที ในขณะที่วิธีการแบบเวิร์คโฟลว์ที่นำเสนอในวิทยานิพนธ์นี้ เป็นการเรียกใช้เว็บเซอร์วิสต่างๆที่กระจายกันอยู่ อินเทอร์เน็ต พบว่าที่จำนวน 2,000 Bootstrapping replicates นั้น เวิร์คโฟลว์ทำงานไม่สำเร็จ เนื่องจากเกิดปัญหาการเกินเวลาของ Data flow ระหว่างเว็บเซอร์วิสดังที่ได้กล่าวแล้ว

สำหรับเวิร์คโฟลว์ใหม่ที่ใช้บริการท้องถิ่นที่พัฒนาขึ้นเอง สามารถแก้ปัญหาการเกินเวลาได้ และสามารถทำงานในการทำนายโครงสร้างต้นไม้สายวิวัฒนาการด้วยจำนวน 2,000 Bootstrapping replicates โดยเวลาเพียง 2.24 นาที ซึ่งเป็นเวลาที่ลดลงอย่างมีนัยสำคัญ และโครงสร้างต้นไม้สายวิวัฒนาการที่ได้มีความถูกต้อง โดยได้รับการตรวจสอบจากศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์

3.6 สรุป

ในบทนี้ได้กล่าวถึงการวาดและการเลือกบริการที่เหมาะสม สำหรับการทำงาน ด้วยวิธีการแบบเวิร์คโฟลว์ โดยให้ข้อเสนอแนะที่เป็นรูปธรรม และใช้เวิร์คโฟลว์ในการวิเคราะห์ สนิปในกุ้ง ที่เป็นงานวิจัยร่วมกับศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์ มหาวิทยาลัยสงขลานครินทร์ เป็นกรณีศึกษา

วิทยานิพนธ์นี้นำเสนอการออกแบบ การพัฒนา และการทดสอบเวิร์คโฟลว์และ เว็บเซอร์วิสสำหรับการวิเคราะห์ สนิป และการสร้างโครงสร้างต้นไม้สายวิวัฒนาการในกุ้งแชบ๊วย ด้วยการใช้เทคโนโลยีมายกริดและทาเวอร์นา บริการท้องถิ่นที่สร้างขึ้นสามารถแก้ไขปัญหาการกิน เวลาและความไม่คงเส้นคงวาในการทำงานได้เป็นอย่างดี และสามารถลดเวลาการทำงานลงได้อย่าง มีนัยสำคัญเมื่อเปรียบเทียบกับการทำงานด้วยวิธีการแบบเดิมหรือคัดลอกแล้ววาง เวิร์คโฟลว์ที่ พัฒนาโดยใช้บริการท้องถิ่นเผยให้เห็นความสัมพันธ์ทางพันธุกรรมของกุ้ง นอกจากนี้ยังสามารถใช้ เวิร์คโฟลว์วิเคราะห์ความสัมพันธ์ทางพันธุกรรมกับสิ่งมีชีวิตอื่นๆได้ ตัวอย่างเช่น สัตว์ทะเล ประเภทหอยและหมีก (Molluse) และประเภทแมลง (Arthropod) โดยผู้ใช้งานสามารถป้อนลำดับ นิวคลีโอไทด์ของสิ่งมีชีวิตเหล่านี้ให้กับเวิร์คโฟลว์ได้ โดยไม่ต้องปรับแต่งเวิร์คโฟลว์ใดๆเพิ่มเติม อีก

บทที่ 4

การออกแบบและการพัฒนาเวิร์คโฟลว์สำหรับการวิเคราะห์สนิปของมนุษย์

4.1 ที่มาและความสำคัญ

การวิเคราะห์หรือทำนายผลกระทบ อันเนื่องมาจากการเปลี่ยนแปลงในระดับระบบนิเวศมนุษย์หรือสนิป ซึ่งกระจายตัวอยู่ประมาณ 10 ล้านตำแหน่งบนโครโมโซม 23 คู่กับการตอบสนองต่อยาหรือภาวะการเป็นโรค จำเป็นต้องใช้ข้อมูลมากมายจากหลากหลายแห่งมาประกอบกันเป็นเวิร์คโฟลว์ และรวบรวมเป็นองค์ความรู้ใหม่เพื่อนำไปใช้ในการศึกษาถึงผลกระทบต่างๆ ของสนิปต่อไป โดยจะใช้เวิร์คโฟลว์นี้เป็นกรณีศึกษาของงานวิทยานิพนธ์

ในบทนี้จะกล่าวถึงประโยชน์ในด้านการพัฒนาเวิร์คโฟลว์ ขั้นตอนการเลือกบริการที่เหมาะสม การออกแบบระบบงานในงานวิจัยด้านสนิปของมนุษย์ การสร้างเวิร์คโฟลว์ของแต่ละกระบวนการ การรวบรวมทุกกระบวนการเป็นเวิร์คโฟลว์เดียวกัน การทดลองวัดเวลาการทำงานของเวิร์คโฟลว์และสรุป

4.2 ประโยชน์ในด้านการพัฒนาเวิร์คโฟลว์

ประโยชน์ในด้านการพัฒนาเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์ ในงานวิจัยเภสัชพันธุศาสตร์มีดังนี้

- รวบรวมข้อมูลทางชีวสารสนเทศที่น่าสนใจ และเป็นประโยชน์จากหลายแหล่งข้อมูล โดยสร้างเป็นเวิร์คโฟลว์ในสิ่งแวดล้อมของโปรแกรมทาเวอร์น่า
- รวบรวมผลการทำงานจากเว็บเซอร์วิสต่างๆ ที่มีความสามารถในการทำนายผลกระทบของการเปลี่ยนแปลงสนิป ที่อาจส่งผลถึงการทำงานของโปรตีนโดยสร้างเป็นเวิร์คโฟลว์ในสิ่งแวดล้อมของโปรแกรมทาเวอร์น่า

- นำเทคโนโลยีด้านชีวสารสนเทศมาประยุกต์ใช้กับการเลือก หรือจัดลำดับ ความสำคัญของสปีที่ควรศึกษา ก่อน และอธิบายผลที่เกิดขึ้นจากความสัมพันธ์ระหว่างสปีกับ กลไกการเป็น โรค หรือการแพ้ชนิดต่างๆ ซึ่งในปัจจุบันยังไม่มีเว็บเซอร์วิสใดเพียงแห่งเดียวที่ สามารถวิเคราะห์หรือทำนายผลกระทบอันเนื่องมาจากการเปลี่ยนแปลงของสปีได้
- ช่วยประหยัดเวลาในการเข้าไปใช้งานเว็บเซอร์วิสต่างๆ ที่ละบริการเพราะ ข้อมูลที่ได้มาจากการศึกษาสปีแบบสแกนหาจีโนมทั้งหมด (Whole genome scan) จะมีปริมาณ มาก หากทำแบบวิธีเก่าหรือคัดลอกแล้ววาง (Copying and Pasting) จะทำได้ยากมาก
- เป็นตัวอย่างของการนำข้อมูลจีโนมมากมายจากหลายแหล่ง มาประยุกต์ใช้ ผ่านกระบวนการชีวสารสนเทศร่วมกับ โปรแกรมทาเวอร์นา และเวิร์ค โฟลว์ที่พัฒนาขึ้นจะเป็นไป เพื่อใช้ประโยชน์ทางการแพทย์ สามารถนำไปสร้างเป็นชุดตรวจกรองการเกิดโรคหรือการแพ้ยา ต่างๆได้ในอนาคต
- ก่อให้เกิดองค์ความรู้ด้านความสัมพันธ์ระหว่างสปีกับการเป็น โรคหรือการ แพ้ชนิดต่างๆ รวมถึงผลกระทบของการเปลี่ยนแปลงสปีที่อาจส่งผลถึงการทำงานของ โปรตีน ซึ่งสามารถนำเวิร์ค โฟลว์ที่ได้สร้างนี้ ไปประยุกต์ใช้กับการวิจัยอื่นๆ ที่เกี่ยวข้องกับสปีจำนวนมาก
- สร้างระบบการรวบรวมข้อมูล ที่ออกแบบเฉพาะสำหรับการวิจัยด้านสปี กับ การเป็น โรคหรือการแพ้ชนิดต่างๆ รวมถึงผลกระทบของการเปลี่ยนแปลงสปีที่อาจส่งผลถึง การทำงานของโปรตีนเป็นระบบแรกของประเทศไทย
- เกิดองค์ความรู้ใหม่ในการใช้เทคโนโลยีชีวสารสนเทศ จากโปรแกรมทาเวอร์- นานในการเชื่อมโยงเว็บเซอร์วิสต่างๆ ที่สำคัญ และมีข้อมูลปริมาณมากเพื่อให้มีความสามารถในการ คำนวณหรือจัดลำดับความสำคัญของสปี และสามารถวิเคราะห์หรือทำนายผลกระทบอัน เนื่องมาจากการเปลี่ยนแปลงของสปี
- เป็นก้าวสำคัญในการประยุกต์ใช้เทคโนโลยีทางคอมพิวเตอร์ มาแก้ปัญหาทาง ด้านวิทยาศาสตร์การแพทย์ ทำให้เกิดความร่วมมือกันระหว่างนักวิจัยด้านวิทยาศาสตร์คอมพิวเตอร์ และวิทยาศาสตร์ชีวภาพ ในที่นี้คือมหาวิทยาลัยสงขลานครินทร์และมหาวิทยาลัยมหิดล

4.3 ขั้นตอนการค้นหาบริการที่เหมาะสม

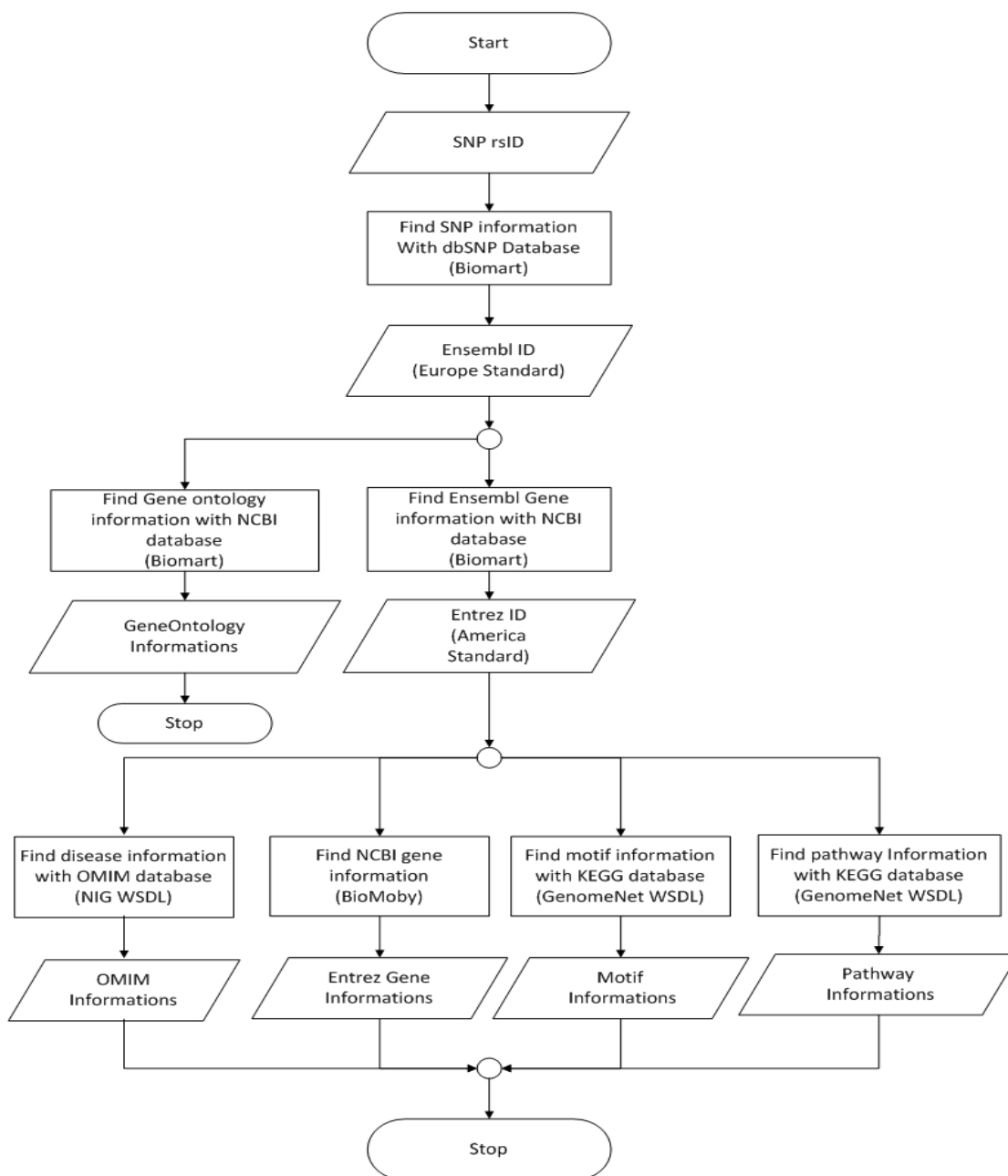
การเลือกบริการที่เหมาะสม ตรงกับความต้องการในการสร้างเวิร์คโฟลว์ได้เสนอไว้แล้วในวิทยานิพนธ์นี้ในบทที่ 3 ซึ่งในบทนี้สามารถใช้วิธีเดียวกันได้

4.4 การออกแบบระบบงานในงานวิจัยด้านสนธิปของมนุษย์

กระบวนการค้นหาข้อมูลในงานวิจัยด้านสนธิปของมนุษย์ ในงานวิทยานิพนธ์นี้เป็นการค้นหาข้อมูลจากฐานข้อมูลต่างๆที่กระจายกันบนอินเทอร์เน็ต และนำไปเป็นพื้นฐานไปสู่การพัฒนาเวิร์คโฟลว์โดยในบทนี้จะกล่าวถึงกระบวนการค้นหาข้อมูล Data flow diagram และการสร้างเวิร์คโฟลว์ดังนี้

4.4.1 กระบวนการค้นหาข้อมูล

การค้นหาข้อมูลของสนธิปในงานวิจัยด้านเกศัชพันธุศาสตร์ สามารถออกแบบการทำงานในภาพรวมในรูปแบบ Flowchart ได้ดังรูปที่ 4.1 ซึ่งมีรายละเอียดดังต่อไปนี้



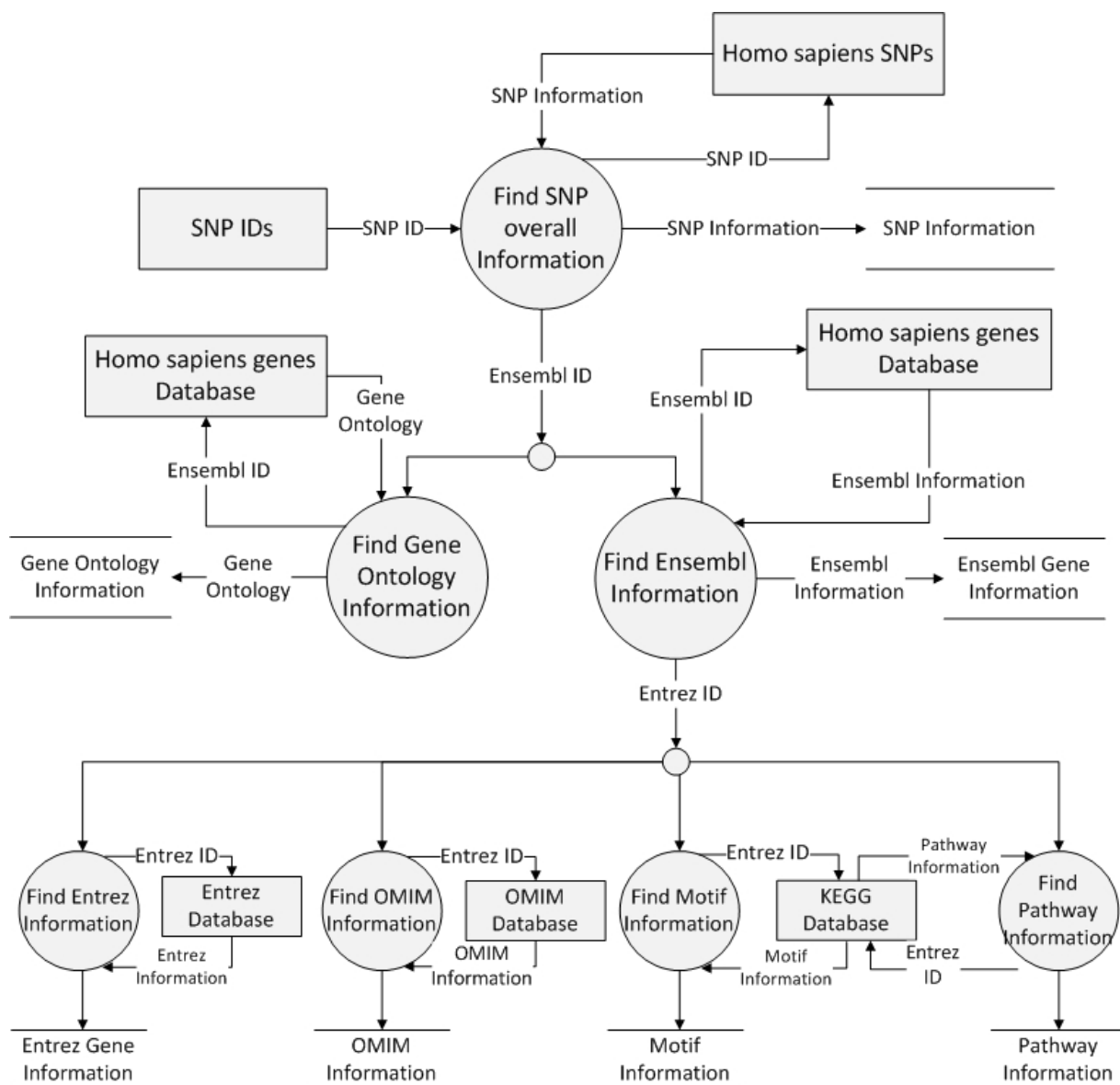
รูปที่ 4.1 Flowchart การทำงานในภาพรวมของเวิร์กโฟลว์สำหรับวิเคราะห์สnpของมนุษย์

1) การค้นหาข้อมูลยีนที่เกี่ยวข้องกับสnpที่สนใจจากฐานข้อมูล dbSNP ของ NCBI โดยใช้ SNP ID จำนวน 500 – 1,000 IDs จากห้องปฏิบัติการจริงซึ่งได้จากการวิเคราะห์ทางหลักสถิติจากจำนวนประมาณ 550,000 ตำแหน่ง ผลลัพธ์ของกระบวนการนี้คือ SNP Information และข้อมูลที่จะนำไปประมวลต่อคือ Ensembl ID [50]

- 2) การค้นหารายละเอียดเกี่ยวกับยีนที่สนใจ จากฐานข้อมูล Homo sapiens genes (NCBI36) ของ NCBI [51][52] เช่น รหัสและชื่อของยีน ผลลัพธ์ของกระบวนการนี้คือ Ensembl Gene Information และข้อมูลที่จะไปประมวลต่อคือ Entrez ID [53]
- 3) การค้นหาข้อมูล Gene Ontology ที่เกี่ยวข้องกับยีนที่สนใจ จากฐานข้อมูล Homo sapiens genes (NCBI36) ของ NCBI ผลลัพธ์ของกระบวนการนี้คือ Gene Ontology Information
- 4) การค้นหาข้อมูลของโรคต่างๆที่เกี่ยวข้องกับยีนที่สนใจ หรือ Online Mendelian Inheritance in Man (OMIM) [54] โดยใช้อินพุตคือ Entrez ID จากฐานข้อมูลของ DDBJ และผลลัพธ์ของกระบวนการนี้คือ OMIM Information
- 5) การค้นหาข้อมูลเกี่ยวกับ Motif [55] ของยีนที่เกี่ยวข้องกับยีนที่สนใจ โดยใช้ อินพุตคือ Entrez ID จากฐานข้อมูล Kyoto Encyclopedia of Gene and Genomes (KEGG) [56] และผลลัพธ์ของกระบวนการนี้คือ Motif Information
- 6) การค้นหาข้อมูลเกี่ยวกับ Pathway หรือกลไกการทำงานของยีนที่สนใจ [57] โดยใช้อินพุตคือ Entrez ID จากฐานข้อมูล KEGG ซึ่งจะแสดงกลไกการทำงานของยีนในรูปแบบของภาพกราฟิกที่เข้าใจง่าย และผลลัพธ์ของกระบวนการนี้คือ Pathway Information
- 7) การค้นหาชื่ออื่นๆของยีนที่สนใจ ที่มีความเกี่ยวข้องหรือใกล้เคียงกันจาก ฐานข้อมูล NCBI และผลลัพธ์ของกระบวนการนี้คือ Entrez Gene Information

4.4.2 Data flow diagram ของเวิร์คโฟลว์

จากกระบวนการค้นหาข้อมูลจากฐานข้อมูลที่ยีน ในที่ต่างๆที่กระจัดกระจายกัน อยู่บนเครือข่ายอินเทอร์เน็ต สามารถอธิบายการไหลและกระบวนการของข้อมูลด้วย Data flow diagram ได้ดังรูปที่ 4.2 ซึ่งจะเห็นว่าอินพุตที่ใช้ในเวิร์คโฟลว์จะเริ่มต้นด้วย SNP IDs ซึ่งผู้ใช้งาน เวิร์คโฟลว์เป็นผู้ระบุ ต่อจากนั้นเวิร์คโฟลว์จะเริ่มทำงานไปตามกระบวนการที่ออกแบบไว้ ซึ่งจะได้ เอาท์พุตสำคัญที่ถูกนำไปใช้งานเป็นอินพุตสำหรับกระบวนการถัดไปคือ จาก SNP IDs ได้เป็น Ensembl IDs และจาก Ensembl IDs ได้เป็น Entrez IDs

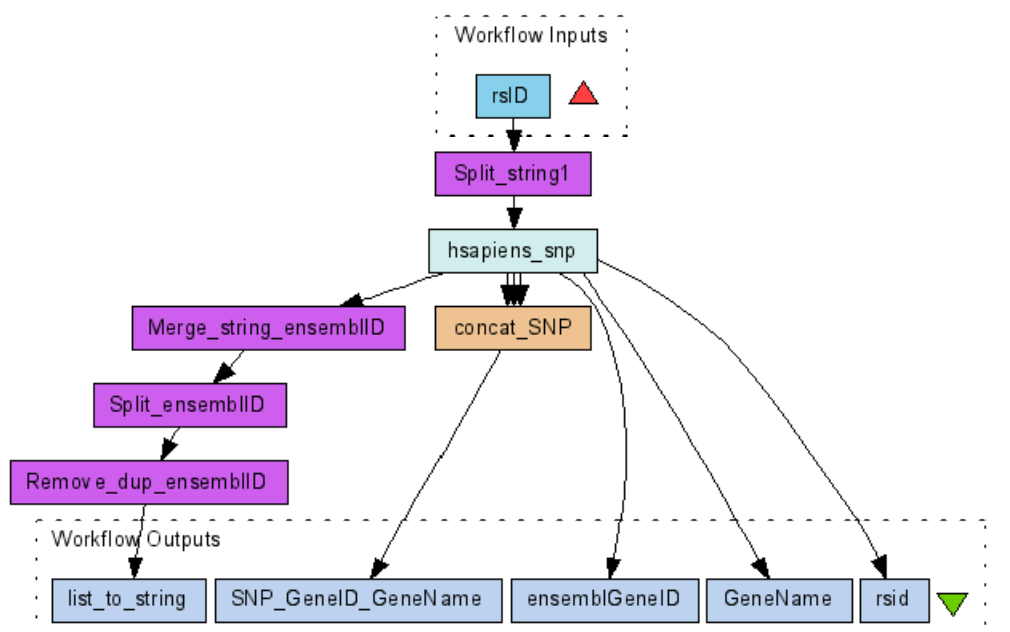


รูปที่ 4.2 Data flow diagram ของเวิร์คโฟลว์

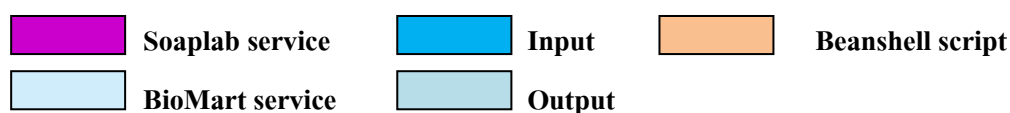
4.5 การสร้างเวิร์คโฟลว์ของแต่ละกระบวนการ (Implementation)

จากการออกแบบ Data flow ดังรูปที่ 4.2 สามารถนำแต่ละกระบวนการมาสร้างเป็นเวิร์คโฟลว์ในงานวิจัยด้านสนิปด้วยโปรแกรมทาวเวอร์น่าได้ดังนี้

1) การค้นหาข้อมูลสลับของมนุษย์ซึ่งใช้ SNP IDs เป็นอินพุตค้นหาจากฐานข้อมูล dbSNP โดยใช้เว็บเซอร์วิสไปโอมาร์ทชื่อ *hsapiens_snp* (Homo sapiens SNPs) เอาท์พุตที่ต้องการคือ SNP Gene ID, Ensembl Gene ID, Gene Name และ rs ID หรือ SNP ID โดยข้อมูลที่จะนำไปเป็นอินพุตให้กับกระบวนการถัดไปคือ Ensembl Gene ID และสร้างเวิร์กโฟลว์ได้ดังรูปที่ 4.3 โดยมีบริการหลักคือ *hsapiens_snp* ทำหน้าที่ค้นหาข้อมูลสลับ แต่ผลลัพธ์การค้นหานั้นยังไม่อยู่ในรูปแบบที่สามารถนำไปใช้เป็นอินพุตกับกระบวนการถัดไปได้ทันที จึงต้องปรับแต่งผลลัพธ์ด้วยบริการท้องถิ่นของ Java (Java local service) เช่น บริการชื่อ *Merge_string_ensemblID* และยังมีการใช้งานบริการท้องถิ่นบีนเชลล์ ชื่อ *concat_SNP* ด้วย ตัวอย่างดังตารางที่ 4.1



Legends:

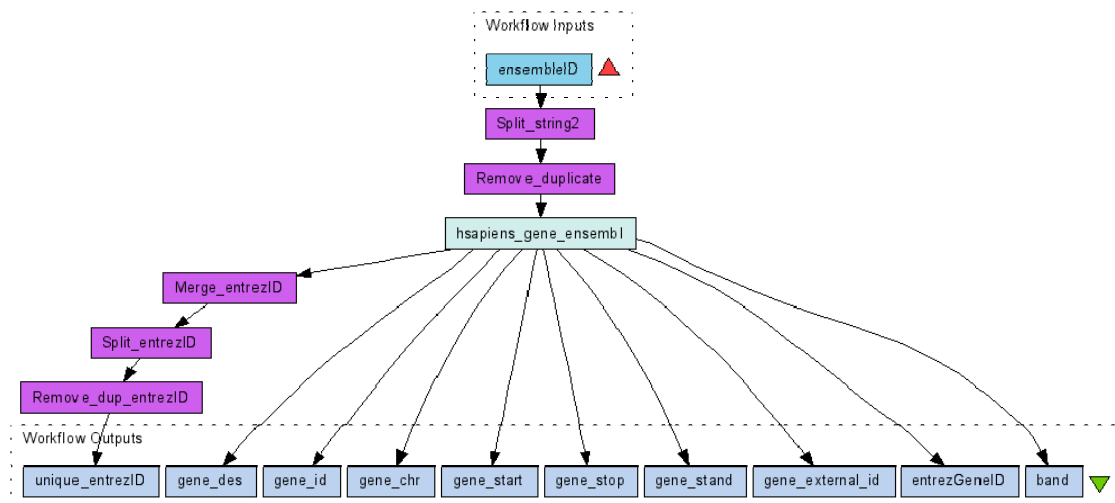


รูปที่ 4.3 เวิร์กโฟลว์การค้นหาข้อมูลสลับที่สนใจ

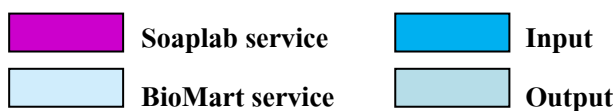
ตารางที่ 4.1 ตัวอย่างการจัดการอินพุตและเอาต์พุตในเวิร์คโฟลว์

ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	ตัวอย่างเอาต์พุต
Split_string1 Split_ensemblID	แยก string	“rs1043502 rs10799658”	“rs1043502” “rs10799658”
Merge_string_ens emblID	รวมหลาย string เป็น string เดียว	“ENSG00000158828” “ENSG00000117242” “ENSG00000117242”	“ENSG00000158828 ENSG00000117242 ENSG00000117242”
Remove_dup_ens emblID	ลบ string ที่ ซ้ำซ้อน	“ENSG00000158828” “ENSG00000117242” “ENSG00000117242”	“ENSG00000158828” “ENSG00000117242”
Concat_SNP	รวบรวม เอาต์พุตเข้า ด้วยกัน	rs1043502 PINK1 ENSG00000158828	“rs1043502,PINK1,ENSG 00000158828”

2) การค้นหาข้อมูลของยีนจากฐานข้อมูล Homo sapiens genes (NCBI36) โดยใช้บริการชื่อ *hsapiens_gene_ensembl* โดยมีอินพุตคือ Ensembl Gene ID ที่ได้เวิร์คโฟลว์ในรูปที่ 4.3 และสร้างเวิร์คโฟลว์ได้ดังรูปที่ 4.4 ซึ่งมีเอาต์พุตดังนี้ Unique Entrez ID, Gene Description, Gene ID, Gene Chromosome, Gene Start Position, Gene Stop Position, Gene Stand, Gene External ID (Gene Name) และ Gene Band. โดยข้อมูลที่จะนำไปเป็นอินพุตให้กับกระบวนการถัดไปคือ Entrez ID มีการใช้บริการท้องถิ่นของ Java เพื่อจัดการกับเอาต์พุตที่จะไปเป็นอินพุตให้กับกระบวนการถัดไป ตัวอย่างดังตารางที่ 4.2



Legends:



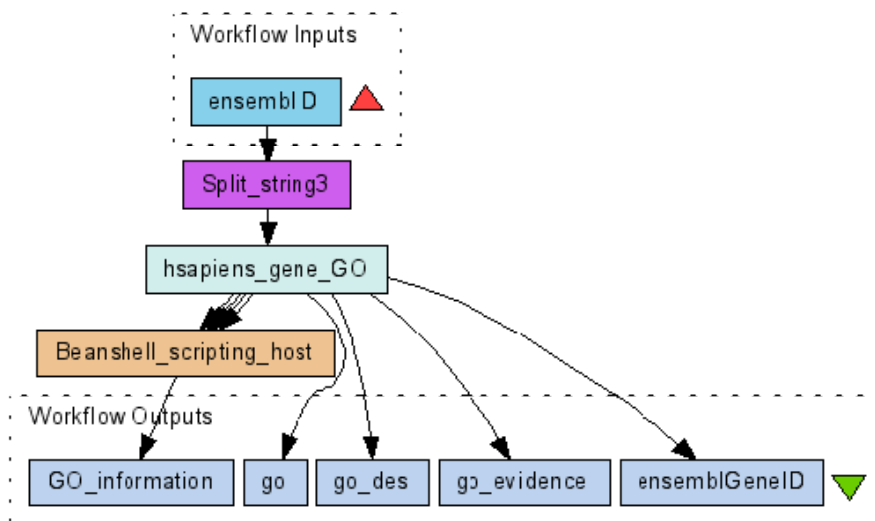
รูปที่ 4.4 เวิร์กโฟลว์การค้นหาข้อมูลของยีน

ตารางที่ 4.2 การทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสท้องถิ่นในการค้นหาข้อมูลยีน

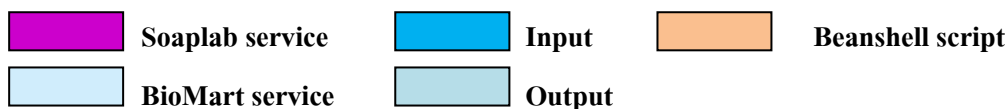
ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	ตัวอย่างเอาต์พุต
Spilt_string2	แยก string	“ENSG00000117242 ENSG00000158828”	“ENSG00000117242” “ENSG00000158828”
Remove_duplicate	ลบ string ที่ซ้ำซ้อน	“ENSG00000117242” “ENSG00000158828”	“ENSG00000117242” “ENSG00000158828”
Hsapiense_gene_ensembl	หาข้อมูลของยีน	“ENSG00000117242” “ENSG00000158828”	“1650” “65018”

3) การค้นหาข้อมูล Gene Ontology จากฐานข้อมูล Homo sapiens genes (NCBI36) โดยใช้บริการไบโอมาร์ทชื่อ *hsapiens_gene_ensembl* ตั้งชื่อเพื่อสื่อความหมายในจุดประสงค์การทำงานว่า *hsapiens_gene_GO* โดยมีอินพุตคือ Ensembl Gene ID ที่ได้จากเวิร์กโฟลว์ในรูปที่ 4.3 และสร้างเวิร์กโฟลว์ได้ดังรูปที่ 4.5 ซึ่งมีเอาต์พุตดังนี้ Gene Ontology, Gene

Description, Gene Ontology Evidence และ Ensembl Gene ID โดยมีบริการท้องถิ่นเป็นเซลล์ รวบรวมเอาที่พูดต่างๆเหล่านี้เข้าด้วยกันเป็นรายงานดังตารางที่ 4.3



Legends:

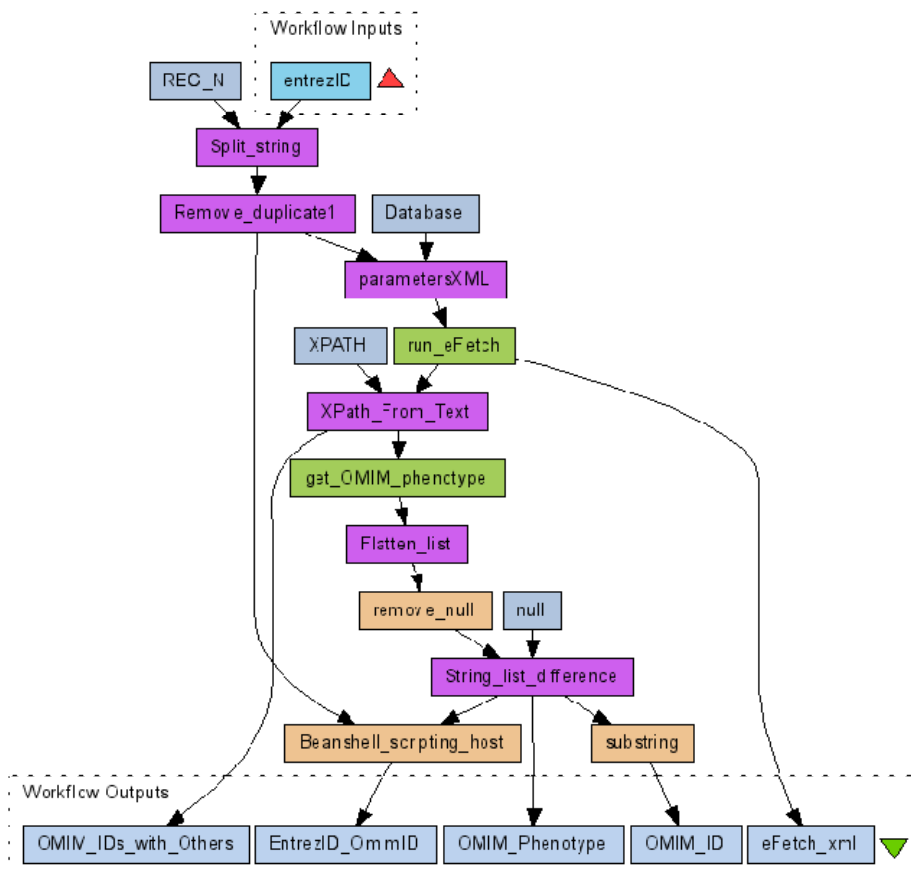


รูปที่ 4.5 เวิร์กโฟลว์การค้นหาข้อมูล Gene Ontology

ตารางที่ 4.3 การทำงานและข้อมูลอินพุตเอาที่พูดของเว็บเซอร์วิสท้องถิ่นในการค้นหาข้อมูล Gene Ontology

ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	ตัวอย่างเอาที่พูด
Spilt_string3	แยก string	“ENSG00000117242 ENSG00000158828”	“ENSG00000117242” “ENSG00000158828”
Beanshell_scriptin g_host	ป็นเซลล์ สคริปต์สำหรับ ต่อ string เข้า ด้วยกัน	gene_id, go, go_des, evi	GO_information = gene_id + "," + go + "," + go_des + "," + evi เช่น ENSG00000117242,GO:0 016020,membrane,IEA

4) การค้นหาข้อมูล OMIM จากฐานข้อมูลของ DDBJ [40] โดยใช้เว็บเซอร์วิสชนิด WSDL ชื่อ *get_OMIM_phenotype* โดยมีอินพุตคือ Entrez Gene ID ที่ได้จากเวิร์คโฟลว์ในรูปที่ 4.3 และสร้างเวิร์คโฟลว์ได้ดังรูปที่ 4.6 โดยมีเอาต์พุตที่สำคัญคือ OMIM ID และ OMIM Phenotype อธิบายตัวอย่างที่สำคัญดังตารางที่ 4.4



Legends:

- | | | |
|---|--|--|
| Soaplab service | Input | Beanshell script |
| Parameter | Output | WSDL service |

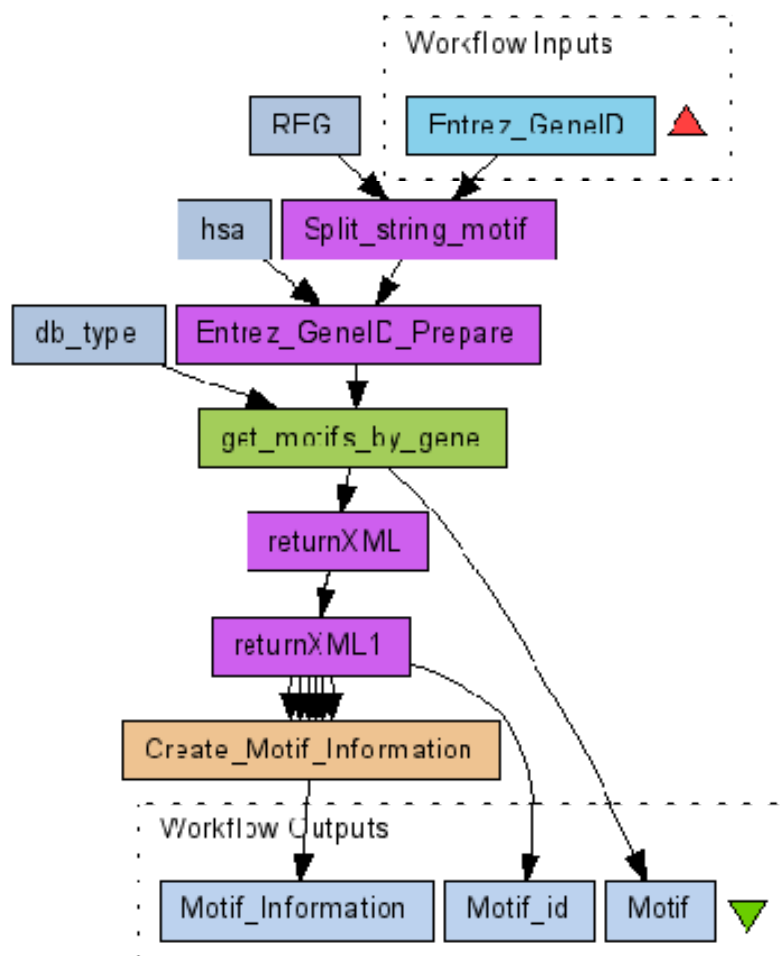
รูปที่ 4.6 เวิร์คโฟลว์การค้นหาข้อมูล OMIM

ตารางที่ 4.4 ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล

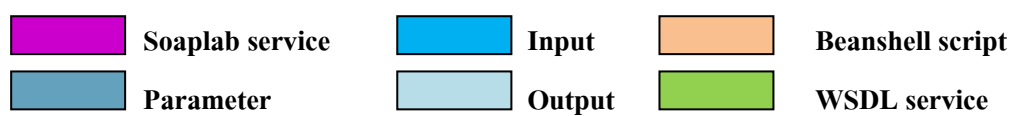
OMIM

ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	ตัวอย่างเอาต์พุต
Spilt_string3	แยก string	“1650 65018”	“1650” “65018”
Run_eFetch	ค้นหาข้อมูลของยีน จากฐานข้อมูล NCBI	“1650” “65018”	ข้อมูลยีนอยู่ในรูปแบบของ ไฟล์ XML ขนาดใหญ่
XPath_From _Text	XPath processor สำหรับดึงอีลิเมนต์ ข้อมูลตามที่ระบุใน Expression	ข้อมูลยีนอยู่ในรูปแบบ ของไฟล์ XML ขนาด ใหญ่ และ XPath Expression ในการดึง OMIM ID	“608309” “605909” “602544”
get_OMIM_ phenotype	ค้นหาข้อมูล OMIM ของยีน โดยใช้ OMIM ID	“608309” “605909” “602544”	OMIM Phenotype แยกเป็น รายการตาม OMIM ID

5) การค้นหาข้อมูล Motif ของยีนจากฐานข้อมูล KEGG โดยใช้เว็บเซอร์วิสชื่อ *get_motifs_by_gene* โดยมีอินพุตคือ Entrez Gene ID ที่ได้จากเวิร์คโฟลว์ในรูปที่ 4.4 และสร้างเวิร์คโฟลว์ได้ดังรูปที่ 4.7 และมีเอาต์พุตคือข้อมูลเกี่ยวกับ Motif ของ Entrez Gene ID ที่สนใจคือ Motif ID, Definition, Gene ID, Start Position, Endposition และ Score ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล Motif ดังตารางที่ 4.5



Legends:

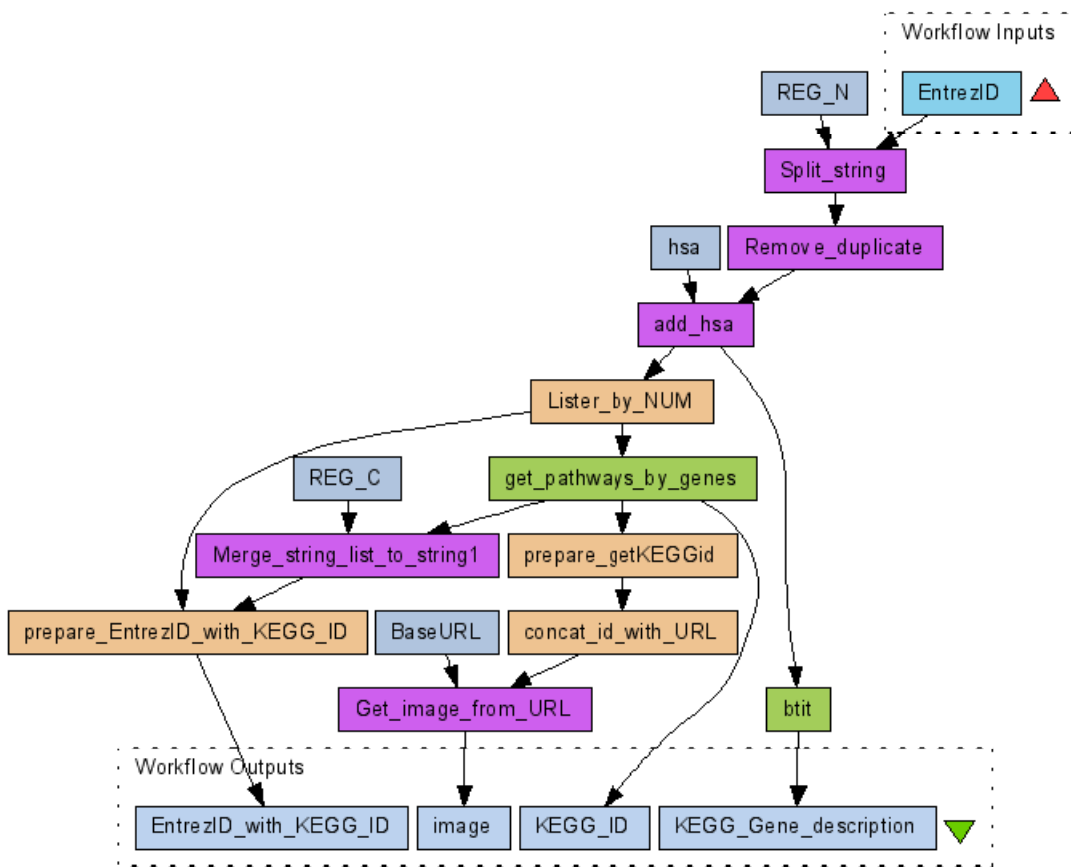


รูปที่ 4.7 เวิร์กโฟลว์การค้นหาข้อมูล Motif ของยีนจากฐานข้อมูล KEGG

ตารางที่ 4.5 ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล Motif

ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	ตัวอย่างเอาต์พุต
Spilt_string_motif	แยก String	“1650 65018”	“1650” “65018”
Entrez_Gene_ID_Prepare	เตรียม Entrez ID ให้เป็น Motif ID โดยการต่อ ‘hsa:’ เข้าไปอยู่ข้างหน้า Entrez ID	“1650” “65018”	“hsa:1650” “hsa:65018”
returnXML returnXML1	ดึงอิลิเมนต์ของข้อมูลต่างๆจากเอกสาร XML ของข้อมูล Motif	เอกสาร XML ของข้อมูล Motif	อิลิเมนต์ข้อมูลคือ motif_id, definition, genes_id, start_position, end_position , score
Create_Motif_Information	บีบเซลล์สคริปต์ในการรวบรวมอิลิเมนต์เข้าด้วยกันเป็นรายงาน	motif_id, definition, genes_id, start_position, end_position , score	motif_info = motif_id + ", " + definition + ", " + genes_id + ", " + start_position + ", " + end_position + ", " + score;

6) การค้นหาข้อมูล Pathways หรือกลไกการทำงานของยีนจากฐานข้อมูล KEGG [8] โดยใช้เว็บเซอร์วิส WSDL ชื่อ *get_pathways_by_genes* โดยมีอินพุตคือ Entrez Gene ID ที่ได้จากเวิร์คโฟลว์ในรูปที่ 4.4 และสร้างเวิร์คโฟลว์ได้ดังรูปที่ 4.8 โดยมีเอาต์พุตที่สำคัญคือ ข้อมูล Pathways และรูปภาพของ Pathways ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล Pathways ดังตารางที่ 4.6



Legends:

- | | | |
|---|--|--|
| Soaplab service | Input | Beanshell script |
| Parameter | Output | WSDL service |

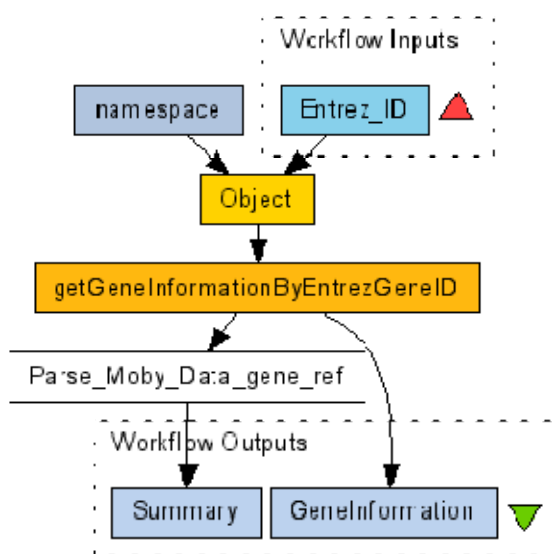
รูปที่ 4.8 เวิร์ดโฟลว์การค้นหา Pathways หรือกลไกการทำงานของยีนจากฐานข้อมูล KEGG

ตารางที่ 4.6 ตัวอย่างการทำงานและข้อมูลอินพุตเอาต์พุตของเว็บเซอร์วิสในการค้นหาข้อมูล

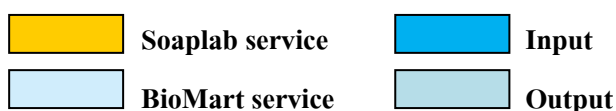
Pathways

ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	ตัวอย่างเอาต์พุต
Spilt_string	แยก string	“1650 65018”	“1650” “65018”
add_hsa	เตรียม Entrez Gene ID ให้เป็น Mofit ID โดยการต่อ ‘hsa:’ เข้าไปอยู่ข้างหน้า Entrez Gene ID	“1650” “65018”	“hsa:1650” “hsa:65018”
Lister_by_N UM	เตรียมข้อมูลอินพุตให้เป็นรายการ (list) สำหรับบริการถัดไป	“hsa:1650” “hsa:65018”	“hsa:1650” “hsa:65018”
get_pathway s_by_genes	ค้นหา Pathways ID จาก Mofit ID	“hsa:7167” “hsa:24”	“path:hsa00010,path:hsa00031,path:hsa00051,path:hsa000710” “path:hsa02010”
Concat_id_ with_URL	เป็นเซลล์สคริปต์สำหรับเชื่อม Pathway ID เข้ากับ Base URL เพื่อดึงข้อมูลภาพกราฟิกของ Pathway	Base URL และ Pathway ID คือ http://www.genome.jp p + "/kegg/pathway/hsa/h sa" + id + ".gif";	เช่น http://www.genome.jp/ kegg/pathway/hsa/hsa000 10.gif
Get_image_f rom_URL	บริการท้องถิ่น Java สำหรับดึงข้อมูลภาพกราฟิกของ Pathway ด้วยคำสั่งขอ GET ของ HTTP	http://www.genome.jp/kegg/pathway/hsa/hsa00010.gif	ภาพกราฟิกของ Pathway ID นั้นๆ

7) การค้นหาข้อมูลของยีนจากฐานข้อมูล Entrez ของ NCBI โดยใช้เว็บเซอร์วิส BioMoby ชื่อ *getGeneInformationByEntrezGeneID* โดยมีอินพุตคือ Entrez Gene ID ที่ได้จากเวิร์คโฟลว์ในรูปที่ 4.4 และสร้างเวิร์คโฟลว์ได้ดังรูปที่ 4.9 ได้เอาท์พุตคือ ข้อมูลเกี่ยวกับ Entrez Gene การใช้บริการ BioMoby จะต้องระบุ namespace ให้กับ Object ของอินพุต ในที่นี้คือ *EntrezGene_EntrezGeneID* เพื่อบอกว่าเป็น Entrez Gene ID และเอาท์พุตที่ได้ก็จะอยู่ในรูปแบบของ Object เช่นเดียวกัน ซึ่งบริการ BioMoby ก็ได้เตรียมตัวดึงข้อมูลเอาท์พุตต่างๆจาก Object เอาไว้ให้แล้วโดยผู้ใช้งานไม่จำเป็นต้องสร้างเว็บเซอร์วิสหรือบีเอ็นเอสสคริปต์เพื่อทำงานดังกล่าว



Legends:



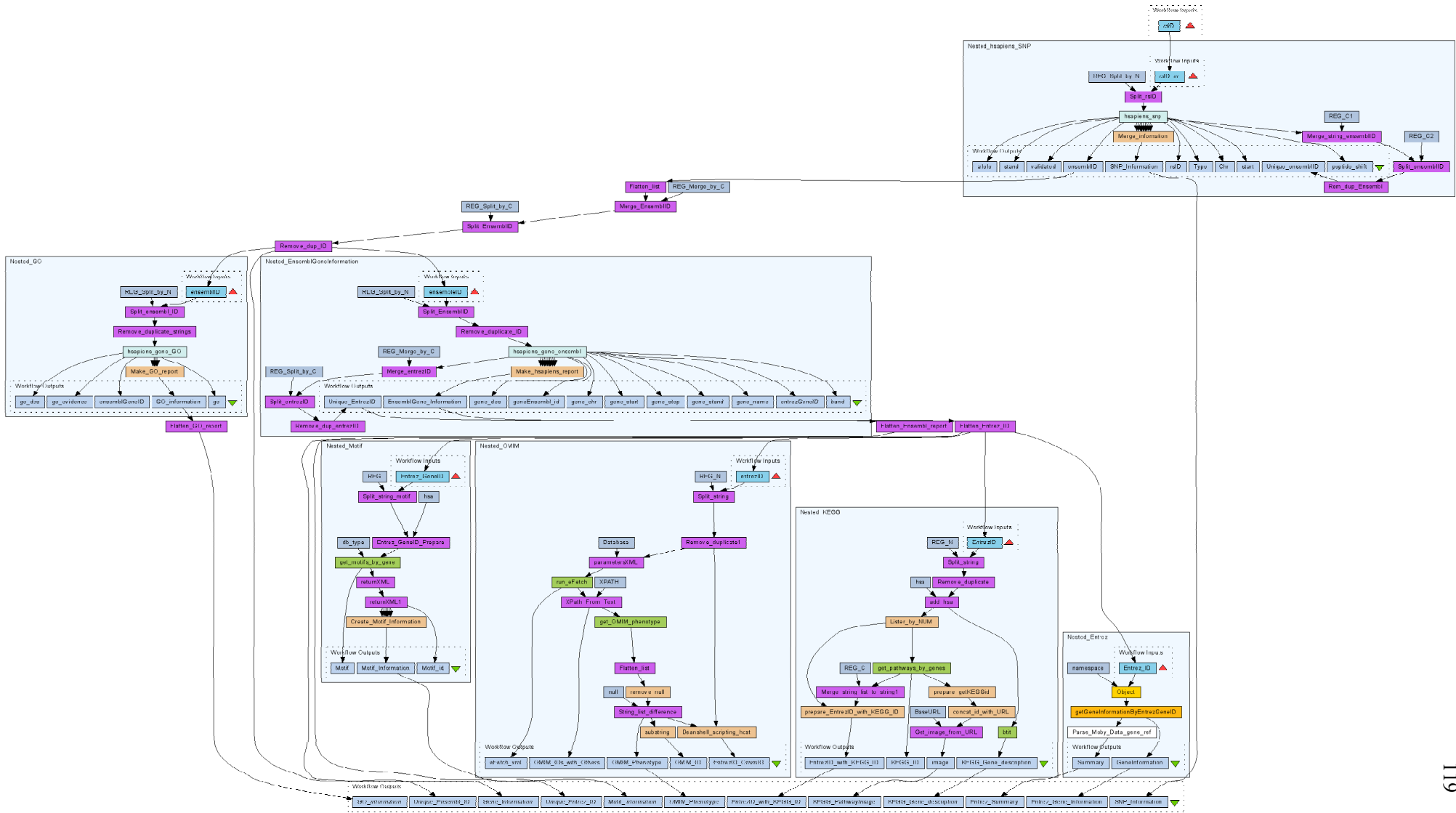
รูปที่ 4.9 เวิร์คโฟลว์การค้นหาข้อมูลของยีนจากฐานข้อมูล Entrez จาก NCBI

4.6 การรวบรวมทุกกระบวนการเป็นเวิร์คโฟลว์เดียวกัน

จากการสร้างและการทดสอบเวิร์คโฟลว์ของแต่ละกระบวนการแล้ว สามารถนำเวิร์คโฟลว์ของแต่ละกระบวนการมารวมไว้ด้วยกัน และปรับแต่งอินพุตและเอาท์พุตตามความเหมาะสม เพราะเวิร์คโฟลว์ในแต่ละกระบวนการเมื่อรวบรวมเข้ามาด้วยกันแล้ว จะมีการทำงานเป็น

ลักษณะของอุปหรือ Nested workflow ดังนั้นเอาต์พุตต่างๆภายใน Nested workflow จึงต้องมีการต่อกรเชื่อมต่อข้อมูล (Data link) ของเอาต์พุตที่สนใจและนำไปใช้งานต่อไป

ในเบื้องต้นได้ทำการเชื่อมต่อทุกเวิร์คโฟลว์ของกระบวนการที่เกี่ยวข้อง ดังที่ได้ทดลองสร้างตามข้อที่ 4.5 ดัง รูปที่ 4.10 โดยมีคุณสมบัติของเวิร์คโฟลว์ดังตารางที่ 4.7



รูปที่ 4.10 เวิร์กโฟลว์การค้นหาข้อมูลสลับ

จากตารางที่ 4.7 แสดงคุณสมบัติของเวิร์คโฟลว์โดยมีบริการทั้งหมด 55 บริการ เป็นบริการไอโอมาร์ทที่ทำหน้าที่คิวรีข้อมูล 3 บริการ คือ *hsapiens_snp*, *hsapiens_ensem* และ *hsapaines_gene_GO* ซึ่งแต่ละบริการเหล่านี้ เป็นส่วนผลิตเอาท์พุตแล้วส่งต่อเอาท์พุตนี้เป็นอินพุต ให้กับบริการอื่นๆอีกหลายบริการ ซึ่งประกอบด้วย 3 ฟิลเตอร์และ 22 แอ็คตริบิวต์และมีดาตาลิงค์ ทั้งหมด 161 ดาตาลิงค์

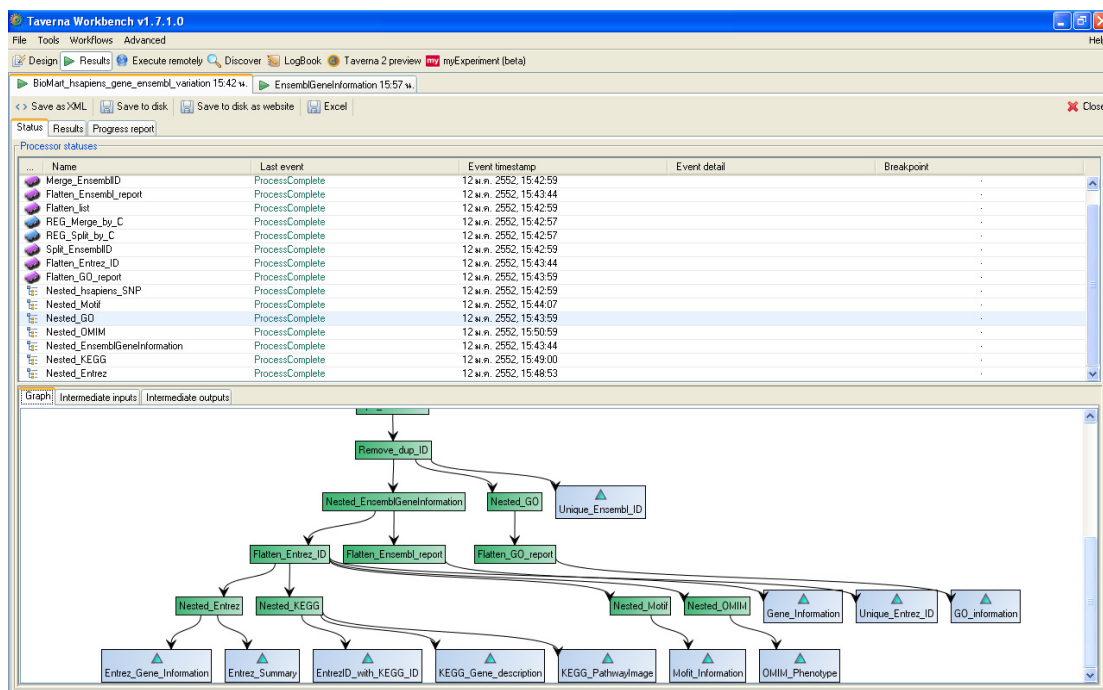
ตารางที่ 4.7 คุณสมบัติของเวิร์คโฟลว์ค้นหาข้อมูลสนิปและระยะเวลาในการทำงาน

รายการ	จำนวน	หมายเหตุ
บริการทั้งหมด	55	-
ดาตาลิงค์ทั้งหมด	161	-
บริการที่เกี่ยวข้องกับบริการไอโอมาร์ท	17	เป็นบริการไอโอมาร์ท 3 บริการ
ดาตาลิงค์ที่เกี่ยวข้องกับบริการไอโอมาร์ท	70	-
ฟिलเตอร์	3	เป็นฟिलเตอร์ของบริการไอโอมาร์ท 3 บริการ
แอ็คตริบิวต์	22	เป็นแอ็คตริบิวต์ของบริการไอโอมาร์ท 3 บริการ
ข้อมูลอินพุต	1,000	SNP IDs
เวลาทำงาน	14.00.50	ชั่วโมง
ผลการทำงาน	ทุกบริการทำงานสำเร็จ	ยกเว้นส่วนการค้นหา Gene Ontology

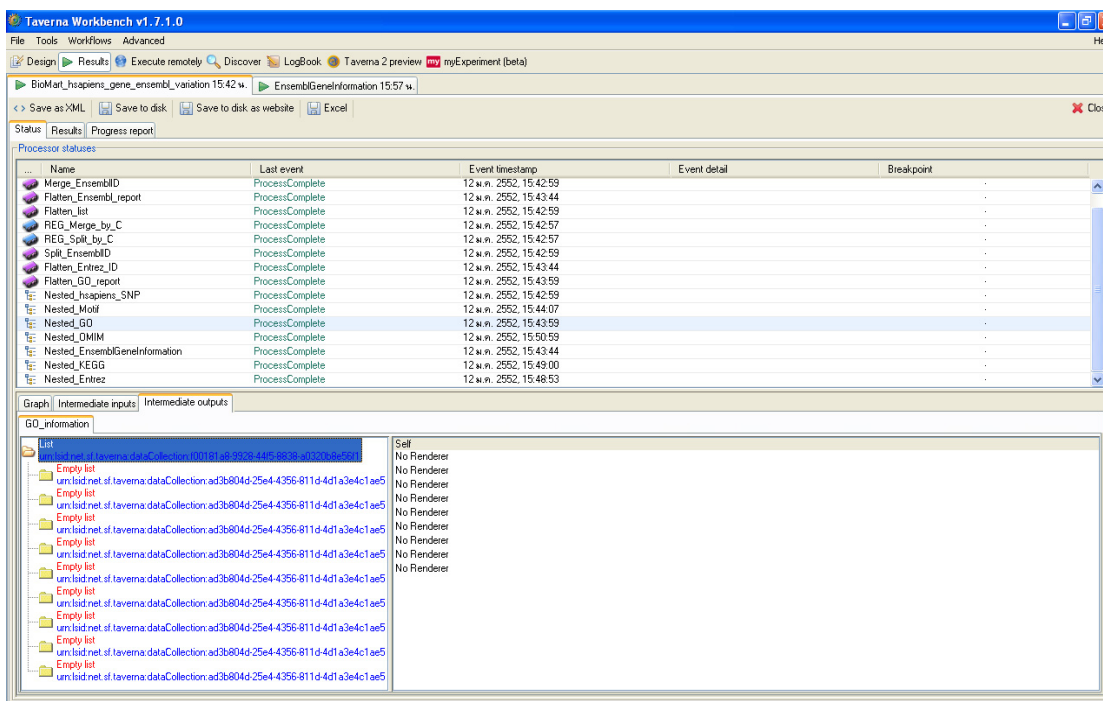
4.7 การทดลองวัดเวลาการทำงานของเวิร์คโฟลว์

การทดลองวัดเวลาการทำงานของเวิร์คโฟลว์โดยใช้ SNP ID จำนวน 1,000 IDs เป็นข้อมูลอินพุต ปรากฏว่าเวิร์คโฟลว์ใช้เวลาประมาณ 14 ชั่วโมง ทุกบริการสามารถทำงานสำเร็จ ยกเว้นในส่วนของการค้นหา Gene ontology ซึ่งโปรแกรมทาเวอร์นาระบุว่าทำงานได้สำเร็จ เช่นเดียวกัน โดยไม่มีรายงานข้อผิดพลาดแต่ไม่ได้ผลิตผลลัพธ์ใดๆ ดังรูปที่ 4.11 และรูปที่ 4.12 ตามลำดับ ซึ่งจริงๆ แล้วคือทำงานไม่สำเร็จ

ในกรณีกระบวนการทำงานตามปกติของเวิร์คโฟลว์นี้ ป้อนอินพุตให้เวิร์คโฟลว์ทำงานเป็น SNP ID จำนวน 1,000 IDs ตามที่ได้ออกแบบไว้ ซึ่งกว่าจะทราบว่าบางส่วนของเวิร์คโฟลว์ทำงานไม่ถูกต้องตรงตามความต้องการก็ต้องรอถึง 14 ชั่วโมง คือรายการที่ไม่มีข้อมูล (Empty list) ก็ต้องเสียเวลารอนกว่าเวิร์คโฟลว์จะทำงานเสร็จ และเมื่อทำงานเสร็จแล้วก็ไม่สามารถที่จะทราบปัญหาได้ทันทีว่าปัญหาเกิดจากสาเหตุอะไร เพราะทุกบริการแจ้งสถานการณ์ทำงานถูกต้อง ซึ่งถัดจากนี้ ก็ต้องตรวจสอบเวิร์คโฟลว์เพื่อหาสาเหตุของข้อผิดพลาด และโดยเฉพาะอย่างยิ่งการผลิตผลลัพธ์เป็นรายการที่ไม่มีข้อมูล (Empty list) จะตรวจสอบได้ยาก เนื่องจากบางกรณีข้อมูลชีวสารสนเทศก็ให้ผลลัพธ์เป็นรายการที่ไม่มีข้อมูลได้ เช่น Ensembl ID หนึ่งๆ ไม่จำเป็นว่าจะต้องมีข้อมูล Entrez ID ด้วยเสมอไป หรือ Entrez ID หนึ่งๆ ก็ไม่จำเป็นว่าจะต้องมีข้อมูล Ensembl ID ด้วยเช่นกัน เนื่องจาก Ensembl เป็นมาตรฐานของยุโรป แต่ Entrez เป็นมาตรฐานของอเมริกา ดังนั้นบางช่วงเวลาระยะหนึ่งๆ ในการปรับปรุงข้อมูลของแต่ละมาตรฐาน จะทำให้ไม่ได้ข้อมูลที่สอดคล้องกันดังกล่าว ซึ่งจะต้องรอนกว่าการจัดการการปรับปรุงนั้นเสร็จสิ้นสมบูรณ์ จึงจะได้ข้อมูลที่ถูกต้องและสอดคล้องเป็นหนึ่งเดียวกัน



รูปที่ 4.11 หน้าต่างแสดงทุกบริการในเวิร์กโฟลว์ทำงานสำเร็จโดยมีแสดงด้วยบล็อกสีเขียว



รูปที่ 4.12 หน้าต่างแสดงส่วนการค้น Gene Ontology ทำงานได้สำเร็จแต่ไม่ได้ให้ผลลัพธ์ใดๆ

4.8 สรุป

ในบทนี้ได้กล่าวถึงการพัฒนาวิธีการทำนายผลกระทบ อันเนื่องมาจากการเปลี่ยนแปลงสนิปต่างๆ รวมทั้งสร้างระบบที่มีความสามารถในการคัดเลือกตำแหน่งสนิปที่น่าสนใจในการนำไปศึกษาต่อทางด้านหน้าที่การทำงานต่อไปด้วยเวิร์คโฟลว์ การแก้ปัญหาทางวิจัยใช้โปรแกรมทาวเวอร์น่าที่มีความสามารถในการสร้างและจัดการรูปแบบข้อมูล รวมทั้งมีความสามารถในการดึงข้อมูลจากแหล่งต่างๆ โดยประสานงานกับเว็บเซอร์วิสชนิดต่างๆและแสดงผลออกมาในรูปแบบที่ง่ายต่อการนำไปศึกษาหรือใช้งานต่อไป และได้รับการตรวจสอบจากนักชีวสารสนเทศแล้วว่าทำงานได้ถูกต้องตรงตามความต้องการและเร็วกว่าการทำงานตามปกติ

แต่จากการทดลองพบว่าบางส่วนของระบบงานในเวิร์คโฟลว์ ยังทำงานได้ไม่ถูกต้องเนื่องจากมีความผิดพลาดมาจากการล้ำสมัยของบริการ ทำให้ผู้ใช้งานต้องใช้เวลายาวนานในการตรวจหาสาเหตุของความผิดพลาดที่แท้จริง ดังนั้นหากเวิร์คโฟลว์ได้รับการตรวจสอบความถูกต้องก่อนการทำงานจริงก็จะสามารถลดเวลาที่อาจจะเสียเป็นชั่วโมงๆหรือเป็นวันได้ โดยในวิทยานิพนธ์นี้จะมุ่งเน้นไปที่บริการไบโอมาร์ทในเวิร์คโฟลว์สำหรับวิเคราะห์ของมนุษย์ เพราะเป็นบริการหลักที่ใช้สำหรับคิวรีข้อมูลซึ่งจะกล่าวในบทต่อไป

บทที่ 5

การออกแบบ การพัฒนาและการทดสอบเวิร์กโฟลว์ สำหรับตรวจสอบความถูกต้องของบริการไอโอมาร์ท

5.1 เกริ่นนำ

โปรแกรมทาเวอร์น่ารุ่นปัจจุบันคือรุ่น 1.7.1 มีปลั๊กอินที่น่าสนใจตัวหนึ่งคือ *Taverna 2 Preview* ซึ่งเป็นการนำคุณสมบัติการทำงานที่เด่นๆของโปรแกรมทาเวอร์น่ารุ่น 2 ที่กำลังพัฒนาอยู่มาให้ทดลองใช้งานก่อน และนอกจากนี้โครงการมายกริดซึ่งเป็นผู้พัฒนาโปรแกรมทาเวอร์น่ายังสร้างกลุ่มผู้ทดสอบ *Taverna 2 Preview* เพื่อแลกเปลี่ยนความเห็น และรับฟังข้อคิดเห็นจากผู้ใช้งานอีกด้วย

จากที่เกริ่นนำไว้ในบทก่อนหน้าว่า ปลั๊กอิน *Taverna 2 Preview* มีฟังก์ชันการทำงาน *Health check report* แต่พบว่าไม่สามารถตรวจสอบความผิดพลาดของเวิร์กโฟลว์ที่สร้างขึ้นเพื่อวิเคราะห์สลิปของมนุษย์ได้ เนื่องจากลักษณะของการออกแบบและการทำงานของเวิร์กโฟลว์มีการทำงานที่ซับซ้อน ประกอบด้วยบริการชนิดต่างๆที่หลากหลายและมีพอร์ตอินพุตและเอาต์พุตจำนวนมาก ทั้งยังมีสถาปัตยกรรมการทำงานแบบโมดูลของเวิร์กโฟลว์ที่ซ้อนกันได้หลายโมดูล (Nested workflow) ทำให้ปลั๊กอิน *Taverna 2 Preview* ไม่สามารถตรวจสอบลงไปถึงโมดูลภายในของเวิร์กโฟลว์ได้ ผู้พัฒนาเวิร์กโฟลว์จึงต้องตรวจสอบเวิร์กโฟลว์เองทีละโมดูลๆโดยใช้ปลั๊กอิน *Taverna 2 Preview* แต่ถึงแม้ว่าจะตรวจสอบทีละโมดูลๆ *Taverna 2 Preview* ก็ยังไม่สามารถตรวจสอบค่าการปรับแต่งต่างๆของบริการไอโอมาร์ท ที่เกิดขึ้นเป็นระยะๆเพื่อให้บริการทันสมัยได้ ซึ่งบริการไอโอมาร์ทเป็นบริการหลักที่เลือกมาใช้ในเวิร์กโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์

เนื่องจากการบริการไอโอมาร์ทมีการปรับปรุงให้ทันสมัยอยู่เสมอ (Upgrade) ทำให้เวิร์กโฟลว์ที่สร้างไว้ก่อนหน้านั้นทำงานได้ไม่ตรงกับความต้องการ คือ ไม่ผลิตผลลัพธ์ใดๆแต่ก็ไม่ได้แจ้งสถานะการทำงานที่ผิดพลาด จากการตรวจสอบพบว่าแท้จริงแล้วปัญหาเกิดจากการเปลี่ยนแปลงของบริการไอโอมาร์ท ซึ่งเป็นการเปลี่ยนแปลงในระดับฟิลเตอร์และแอ็คตริบิวต์ ก็ต้องใช้เวลาในการตรวจสอบ คืออาจจะเป็นวันๆหรือหลายชั่วโมง ยิ่งกว่านั้นผู้วิจัยพบว่า เมื่อรู้

แล้วว่าบริการไปโอมาร์ทใดที่ล้ำสมัยแล้ว ไม่ว่าผู้ใช้งานจะพยายามปรับแต่งวิธีการเรียกใช้บริการ โอโม่มาร์ทอย่างไรก็ตาม ก็ไม่สามารถแก้ปัญหาได้ ทำให้ยิ่งเสียเวลาไปกับการปรับแต่งอีก เพราะ ยังไม่มีองค์ความรู้หรือคำแนะนำวิธีแก้ไขปัญหาที่ถูกต้อง

เนื้อหาของบทนี้จะกล่าวถึงแนวคิดในการออกแบบ การพัฒนาและการทดสอบ เวิร์คโฟลว์ เพื่อตรวจสอบการเปลี่ยนแปลงของบริการไปโอมาร์ทในเวิร์คโฟลว์สำหรับวิเคราะห์ สนิปแบบอัตโนมัติ โดยจะสามารถระบุได้ว่าบริการและฟิลเตอร์หรือแอ็คตริบิวต์ตัวใดที่ล้ำสมัยไป แล้วได้ และนั่นย่อมหมายถึงบริการไปโอมาร์ทนั้นได้ล้ำสมัยไปแล้วเช่นกัน

5.2 แนวคิดการออกแบบเวิร์คโฟลว์สำหรับการตรวจสอบบริการไปโอมาร์ท

ในส่วนนี้ จะอธิบายถึงแนวคิดการออกแบบเวิร์คโฟลว์สำหรับการตรวจสอบ บริการไปโอมาร์ทโดยจะให้รายละเอียดโครงสร้างของ Scufi ซึ่งเป็นเอกสารที่ใช้บันทึกเวิร์คโฟลว์ และโครงสร้างของบริการไปโอมาร์ทที่ใช้ในงานวิจัย ดังนี้

5.2.1 โครงสร้างของ Scufi

Scufi ย่อมาจาก Simple conceptual unified flow language เป็นภาษาและมาตรฐาน ของโปรแกรมทาเวอร์น่าที่มายกริดสร้างขึ้นมา [9] เอกสาร Scufi อธิบายอิลิเมนต์และคุณสมบัติที่ ปรากฏในไฟล์เวิร์คโฟลว์โดยใช้ไวยากรณ์ XScufi XML ไวยากรณ์นี้ใช้ในโปรแกรม ทาเวอร์น่าใน การบันทึกหรือเรียกใช้การประกาศโครงสร้างของเวิร์คโฟลว์ โครงการมายกริดไม่สนับสนุนการใช้ การภาษานี้ นอกเหนือจากสิ่งแวดล้อมของโปรแกรมทาเวอร์น่า

- **Scufi top level tag**

Tag ระดับบนสุดสำหรับการประกาศ Scufi หนึ่งๆ ประกอบด้วยอิลิเมนต์ ต่างๆที่เป็นสมาชิกในเวิร์คโฟลว์คือ บริการ (Processor), การเชื่อมโยง (Link), แหล่งผลิตข้อมูล (Source), แหล่งใช้ข้อมูล (Sink) และการประสานงาน (Coordination) สำหรับงานวิจัยนี้จะเป็นการวิเคราะห์อิลิเมนต์ที่อยู่ใน Tag

processor เท่านั้น Tag ต่างๆเหล่านี้จะถูกอธิบายภายใต้เนมสเปซ (Namespace) <http://org.embl.ebi.escience/xscufl/0.1alpha> โดยมีโครงสร้างดังนี้

- `<scufl version(string) log(int)>` ประกาศว่าเป็นเอกสาร Scufi และระบุรุ่นของภาษา
 - `<processor>*` ประกาศกระบวนการทำงานแบบอะตอมมิกในเวิร์ลด์โพลว์
 - `<link>*` ประกาศการเชื่อมโยงของพอร์ตข้อมูลอินพุตและเอาต์พุตต่างๆในเวิร์ลด์โพลว์
 - `<source>*` ประกาศพอร์ตข้อมูลที่ทำหน้าที่เป็นอินพุต
 - `<sink>*` ประกาศพอร์ตข้อมูลที่ทำหน้าที่เป็นเอาต์พุต
 - `<coordination>*` ประกาศการเชื่อมโยงแบบควบคุมของบริการภายในเวิร์ลด์โพลว์

● โหนดบริการ

การตรวจสอบบริการไป โอมาท์ที่เสนอในวิทยานิพนธ์ มีข้อมูลที่เกี่ยวข้องคือ Processor name หรือชื่อบริการ, Dataset name หรือชื่อดาตาเซต, Filter name หรือชื่อฟิลเตอร์และ Attribute name หรือชื่อแอตทริบิวต์ซึ่งอิลิเมนต์ของข้อมูลเหล่านี้ เป็นสมาชิกของโหนดบริการ(Processor node) ทั้งสิ้น โครงสร้างของโหนดบริการคือ

- `<processor name(string)>` ชื่อบริการ โดยคำปรียายชื่อบริการจะเป็นชื่อเดียวกับดาตาเซต และผู้ใช้งานสามารถเปลี่ยนชื่อให้สื่อความหมายตามการทำงานหรือตามความต้องการได้
 - `<description>?`
 - textual description (PCDATA) คำอธิบายการทำงานของบริการ
 - `<SPEC ELEMENT maxretries(int)? retrydelay(int)? retrybackoff(double)?>` หน่วยการทำงานของบริการ ได้แก่ ค่าคงที่, บริการแบบ WSDL, บริการท้องถิ่น, บริการ Soaplab, บริการ

BioMoby หรือบริการ BioMart และรวมถึง โมดูลของเวิร์คโฟลว์ที่ ซ้อนได้หลายโมดูล (Nested workflow) และคุณสมบัติการทำงานซ้ำ ของบริการเหล่านี้หรือ Retry ค่าโดยปริยายคือ maxretries และ retrydelay = 0 ส่วน retrybackoff = 1

- <iterationstrategy>? กลไกการทำซ้ำของพอร์ตข้อมูลสำหรับการ ป้อนเป็นอินพุตมีสองชนิดคือ cross product (หลายต่อหลาย) และ Dot product (หนึ่งต่อหนึ่ง) โดยมี Namespace คือ <http://org.embl.ebi.escience/xscufliteration/0.1beta10>

เวิร์คโฟลว์ที่สร้างด้วยโปรแกรมทาวเวอร์นั้น จะบันทึกเป็นไฟล์เอกสาร Scufi ซึ่งเป็นมาตรฐานที่มายกริดสร้างขึ้น อาจกล่าวได้ว่า Scufi เป็นเอกสาร XML ชนิดหนึ่งก็ได้ เนื่องจากมีโครงสร้างแบบภาษา XML จึงสามารถใช้ XML Path Language (XPath) [58] เข้าไปถึง ข้อมูลที่ต้องการต่างๆของบริการไบโอมาร์ทที่อยู่ในเอกสาร Scufi ได้ โปรแกรมทาวเวอร์เองก็มี บริการท้องถิ่นคือ *'XPath From Text'* ไว้ให้ทำงานดังกล่าวได้ดังรูปที่ 5.1 ซึ่งสามารถใส่อินพุตเป็น เอกสาร XML และ XPath expression เพื่อดึง (Extract) ข้อมูลจากเอกสาร XML อินพุตนั้นได้เลย โดยในวิทยานิพนธ์นี้ เอกสาร XML ดังกล่าวก็คือเอกสารเวิร์คโฟลว์หรือเอกสาร Scufi ของเวิร์ค-โฟลว์

XPath_From_Text	
xpath	odelist
xml-text	odelistAsXML

รูปที่ 5.1 บริการท้องถิ่น *'XPath From Text'* ในโปรแกรมทาวเวอร์

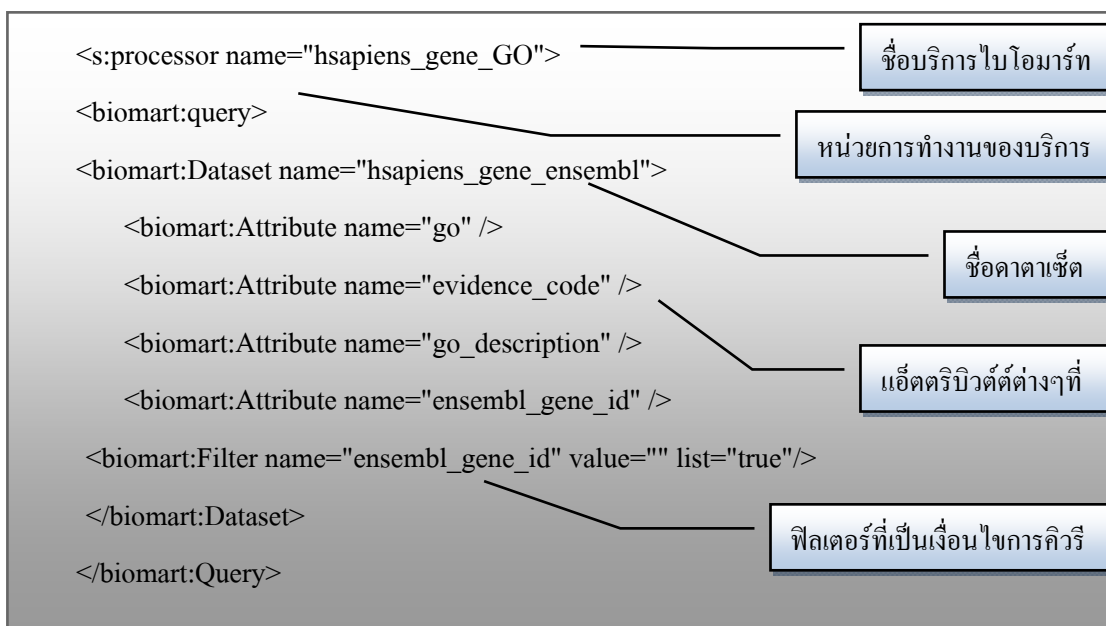
5.2.2 โครงสร้างของบริการไบโอมาร์ทที่ใช้งานวิจัย

บริการไบโอมาร์ทหนึ่งๆจะประกอบด้วยหลายคอมโพเนนท์ หนึ่งในนั้นคือชื่อของบริการหรือ Processor/Service name เนื่องจากบริการไบโอมาร์ททำงานในลักษณะคล้ายเหมืองข้อมูล (Data mining) [12] ดังนั้น บริการไบโอมาร์ทสามารถคิวรีข้อมูล จากฐานข้อมูลหรือดาตาเซตที่ผู้ให้บริการกำหนดไว้ได้ โดยในแต่ละดาตาเซตก็จะประกอบด้วยฟิลเตอร์ (ฟิลด์ข้อมูลของเงื่อนไขในการคิวรี) และแอตทริบิวต์ (ฟิลด์ข้อมูลของเอาต์พุตจากการคิวรี) ที่แตกต่างกันออกไป งานวิจัยนี้ได้ตีกรอบเนื้อหาการตรวจสอบของบริการไบโอมาร์ทไว้ว่า บริการหนึ่งๆจะมีหนึ่งชื่อบริการ และสามารถคิวรีข้อมูลได้หนึ่งดาตาเซต เพื่อให้เป็นไปตามแนวคิดการทำงานแบบอะตอมมิกหรือปรมาณูของเว็บเซอร์วิส และจะง่ายต่อการดีบั๊กหรือตรวจสอบการทำงาน และดาตาเซตหนึ่งๆประกอบด้วยหลายฟิลเตอร์และหลายแอตทริบิวต์ ดังแสดงโครงสร้างบริการไบโอมาร์ทดังรูปที่ 5.2

BioMart Service	
A Processor/Service Name	
A Dataset	
Filters (Query fields)	Attributes (Output fields)

รูปที่ 5.2 โครงสร้างของบริการไบโอมาร์ทในสิ่งแวดล้อมของทาวเวอร์นา

ตัวอย่างเอกสาร Scuff ที่มีการบันทึกอิลิเมนต์ของบริการไบโอมาร์ทเช่น ชื่อบริการคือ *hsapiens_gene_GO* โดยคิวรีที่ดาตาเซตชื่อ *hsapiens_gene_ensembl* และมีฟิลเตอร์คือ *ensembl_gene_id* และแอตทริบิวต์ต่างๆคือ *go*, *evidence_code*, *go_description*, และ *ensemble_gene_id* ดังรูปที่ 5.3



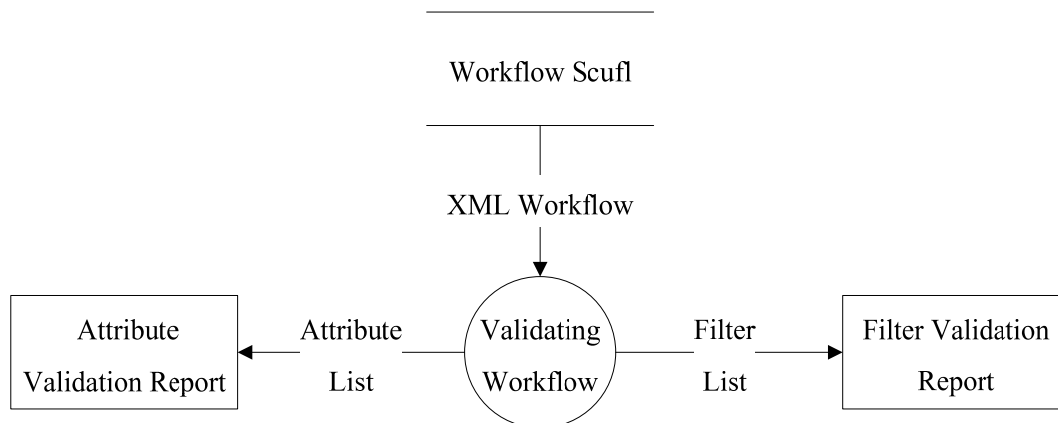
รูปที่ 5.3 ตัวอย่างเอกสาร Scufi ที่มีการบันทึกอิลิเมนต์ของบริการไบโอมาร์ท

อนึ่ง จากกลไกการทำงานของโปรแกรมทาวเวอร์น่าเอง เมื่อเปิดเวิร์คโฟลว์ใดขึ้นมา โปรแกรมจะทำการตรวจสอบการตอบกลับของเอนพอยต์โดยอัตโนมัติ (Responding) ซึ่งหากเอนพอยต์ของบริการใด ไม่มีการตอบกลับ โปรแกรมทาวเวอร์น่าจะแจ้งเตือน และไม่สามารถเปิดเวิร์คโฟลว์นั้นขึ้นมาให้ทำงานได้โดยปริยาย [8][9] ดังนั้นในการออกแบบการตรวจสอบบริการไบโอมาร์ท ในวิทยานิพนธ์นี้จึงไม่ได้สนใจการตรวจสอบในเรื่องนี้อีก เพราะถือว่าโปรแกรมทาวเวอร์น่าได้ทำหน้าที่ตรวจสอบเรื่องนี้อยู่แล้ว แต่จะเน้นความสนใจไปที่ การตรวจสอบฟิลเตอร์และแอตทริบิวต์ของคาตาเซตที่เกี่ยวข้องกับบริการไบโอมาร์ทนั้นๆ

5.3 การออกแบบ Data flow ของเวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท

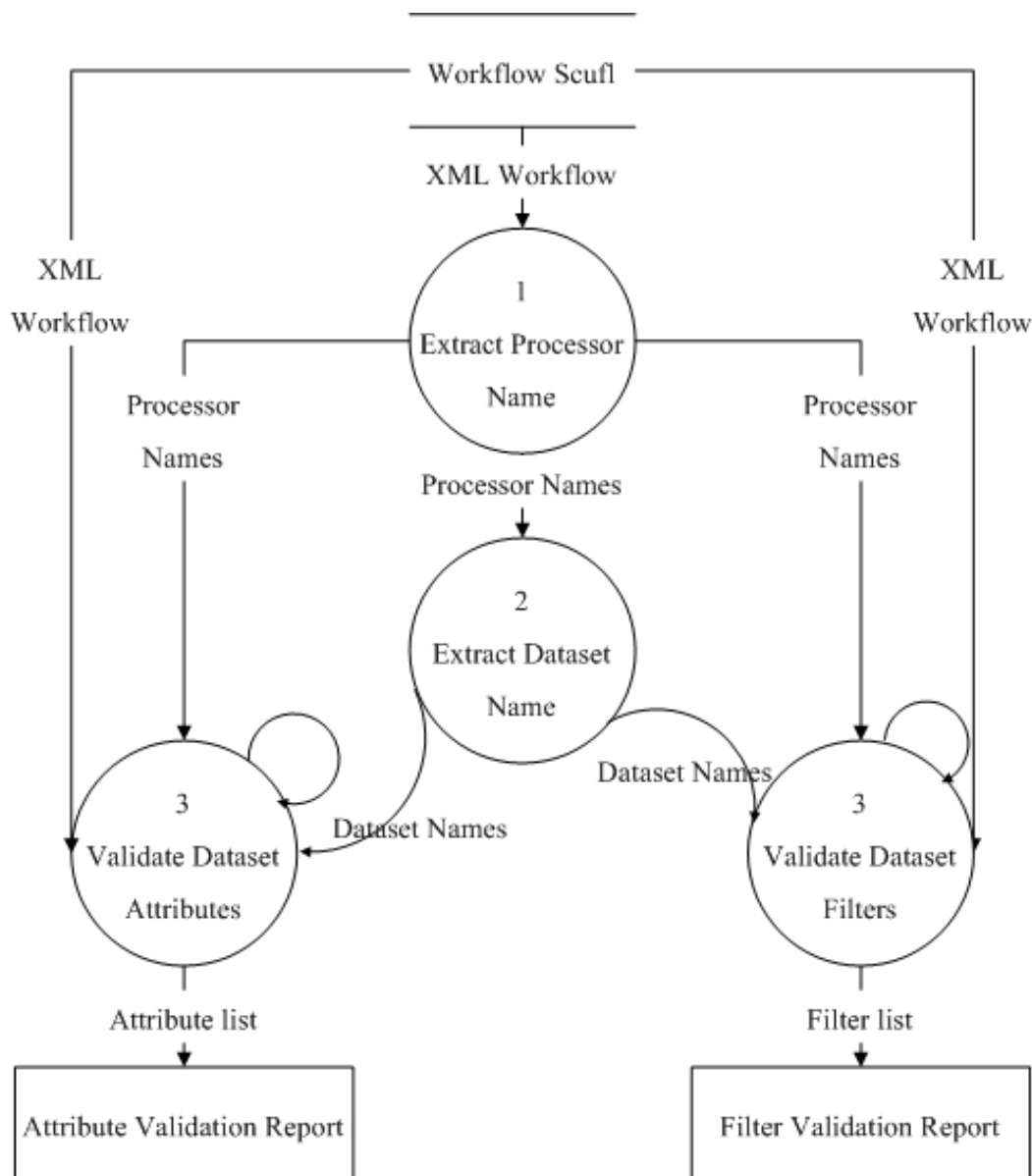
ในขั้นนี้จะอธิบายการออกแบบการทำงานของเวิร์คโฟลว์ที่จะใช้ในการตรวจสอบบริการไบโอมาร์ทในรูปแบบ Data flow diagram เพื่อให้เห็นการไหลของข้อมูลภาพในเวิร์คโฟลว์ประกอบด้วย การทำงานของเวิร์คโฟลว์ในภาพรวมหรือ Context diagram และกระบวนการตรวจสอบบริการไบโอมาร์ทดังนี้

5.3.1 การทำงานของเวิร์คโฟลว์ในภาพรวม



รูปที่ 5.4 Context diagram ของระบบการตรวจสอบของเวิร์คโฟลว์

จากรูปที่ 5.4 แสดง Context diagram ของระบบการตรวจสอบเวิร์คโฟลว์ ข้อมูลอินพุตคือเอกสาร ScufI ของเวิร์คโฟลว์ที่มีการเรียกใช้งานบริการไปโอมาร์ท โดยอาจจะอยู่ในรูปแบบของเวิร์คโฟลว์หรือโมดูลของเวิร์คโฟลว์ที่ซ้ำซ้อนกันหลายโมดูล (Nested workflow) ก็ได้ เพื่อนำเข้าสู่กระบวนการตรวจสอบเวิร์คโฟลว์ แล้วจะแจ้งผลลัพธ์เป็นรายงานการตรวจสอบบริการไปโอมาร์ทเพื่อให้ผู้ใช้งานนำข้อมูลไปตัดสินใจเกี่ยวกับการทำงานต่อไปได้ โดยกระบวนการการทำงานของการตรวจสอบจะอธิบายในหัวข้อถัดไป

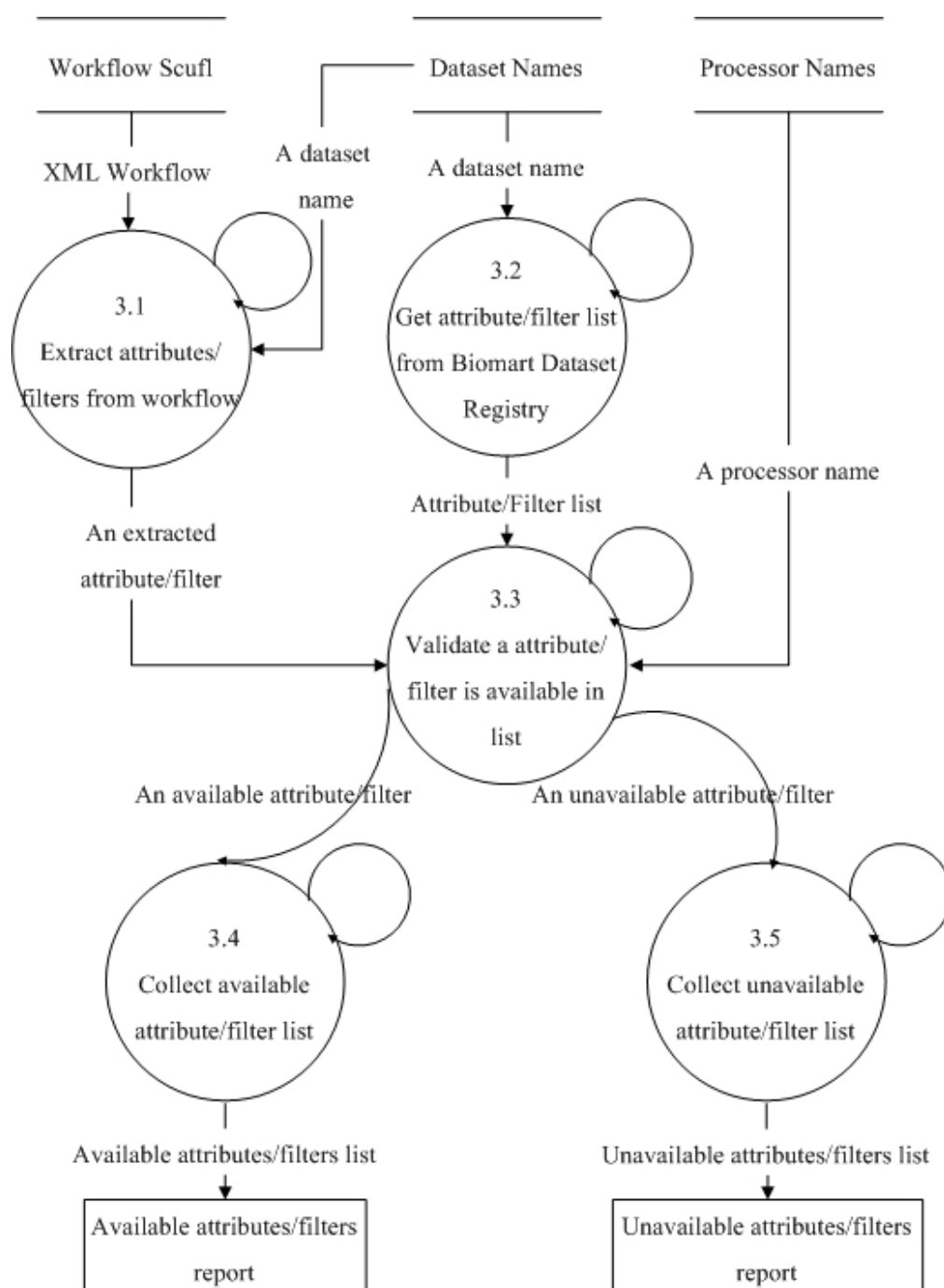


รูปที่ 5.5 Data flow diagram ชั้นที่ 2 อธิบายกระบวนการตรวจสอบบริการไปโอมาร์ท

จากรูปที่ 5.5 แสดงขั้นตอนและการไหลของข้อมูลในการตรวจสอบบริการไปโอ-มาร์ทโดยมีอินพุตคือเอกสาร Scuf ของเวิร์คโฟลว์ จากนั้นทำการดึง (Extract) ชื่อของบริการไปโอ-มาร์ททุกตัวที่อยู่ในเวิร์คโฟลว์ เรียกว่า Processor name แล้วใช้ Processor name ไปหาชื่อดาตาเซตอีกครั้งหนึ่ง เมื่อได้ชื่อดาตาเซตของบริการนั้นๆแล้ว ก็ใช้ชื่อดาตาเซตเป็นคำสำคัญเพื่อเข้าไปตรวจสอบทั้งฟิลเตอร์และแอ็คตริบิวต์ เพื่อหาว่าข้อมูลใดที่ยังทันสมัยอยู่ และข้อมูลใดที่ล้าสมัยแล้ว

กลไกการตรวจสอบฟิลเตอร์และแอตทริบิวต์ที่มีการทำงานแบบวนซ้ำหรือลูป สามารถอธิบายได้ดังรูปที่ 5.6

5.3.2 กลไกการตรวจสอบบริการไบโอมาร์ท



รูปที่ 5.6 Data flow diagram ขั้นที่ 3 อธิบายกระบวนการตรวจสอบบริการไบโอมาร์ท

จากรูปที่ 5.6 เป็นกลไกในการทำหน้าตรวจสอบบริการไปโอมาร์ท โดยมีอินพุต 3 ชนิดคือ เอกสาร Scufi ของเวิร์คโฟลว์, ชื่อบริการไปโอมาร์ททั้งหมด และชื่อคาตาเซ็ทของบริการนั้นๆ การทำงานจะเริ่มด้วยการส่งข้อมูลเข้าไปที่ละบริการหรือที่ละคาตาเซ็ท เพราะหนึ่งบริการจะมีเพียงหนึ่งคาตาเซ็ทเท่านั้น ดังนั้นการทำงานจึงมีลักษณะเป็นลูปหรือวนทำงานซ้ำๆจนกว่าข้อมูลอินพุตจะหมด โดยจะใช้ชื่อคาตาเซ็ทเข้าไปดึงอิลิเมนต์ของฟิลเตอร์และแอ็คตริบิวต์ที่อยู่ในเอกสาร Scufi ของเวิร์คโฟลว์ที่ต้องการจะตรวจสอบ ในอีกด้านหนึ่งก็ใช้ชื่อคาตาเซ็ทอีกเช่นกัน ในการไปสืบค้นข้อมูลรายชื่อฟิลเตอร์และแอ็คตริบิวต์ ที่ทันสมัยจากฐานข้อมูลริจิสทรีของไปโอมาร์ทที่ให้บริการบนเครือข่ายอินเทอร์เน็ต

เมื่อได้ข้อมูลทั้ง 2 ส่วนแล้ว จึงเปรียบเทียบว่า ฟิลเตอร์และแอ็คตริบิวต์ใดๆ จากเอกสาร Scufi ของเวิร์คโฟลว์นั้น ปรากฏอยู่ในรายชื่อฟิลเตอร์และแอ็คตริบิวต์ของริจิสทรีไปโอมาร์ทหรือไม่ หากปรากฏก็แสดงว่า ฟิลเตอร์และแอ็คตริบิวต์นั้นๆยังคงทันสมัยอยู่ ก็จะรวบรวมข้อมูลเกี่ยวกับฟิลเตอร์และแอ็คตริบิวต์นี้ ได้แก่ คำอธิบายและรายละเอียดต่างๆ นำเสนอต่อผู้ใช้งานเพื่อเก็บไว้เป็นข้อมูลอ้างอิงสำหรับการตรวจสอบในอนาคต เพราะเมื่อเวลาผ่านไป บริการนั้นๆก็อาจจะล้าสมัยได้อีกเช่นกัน การนำเสนอข้อมูลคำอธิบายฟิลเตอร์และแอ็คตริบิวต์นี้ ก็จะช่วยให้สามารถสืบค้นหรือเปรียบเทียบได้ในอนาคต

ส่วนฟิลเตอร์และแอ็คตริบิวต์ใด ที่ไม่ปรากฏในอยู่รายชื่อฟิลเตอร์และแอ็คตริบิวต์ของริจิสทรีไปโอมาร์ท แสดงว่า บริการไปโอมาร์ทที่มีฟิลเตอร์และแอ็คตริบิวต์เหล่านั้นได้ล้าสมัยไปแล้ว จึงไม่มีข้อมูลใดๆที่เกี่ยวข้องกับฟิลเตอร์และแอ็คตริบิวต์นั้นๆอีก จึงรายงานเพื่อแจ้งให้ผู้ใช้งานทราบ

5.4 การพัฒนาเวิร์คโฟลว์

ในขั้นนี้จะอธิบายขั้นตอนการพัฒนาเวิร์คโฟลว์ (Workflow implementation) โดยจะให้รายละเอียดการดึงอิลิเมนต์ข้อมูลต่างๆในเอกสาร Scufi ของเวิร์คโฟลว์ การสืบค้นข้อมูลรายชื่อฟิลเตอร์และแอ็คตริบิวต์ที่ทันสมัย การเปรียบเทียบข้อมูลในเวิร์คโฟลว์กับริจิสทรีของไปโอมาร์ท และการปรับแต่งกลไกการซ้ำของพอร์ตอินพุตในเวิร์คโฟลว์ดังต่อไปนี้

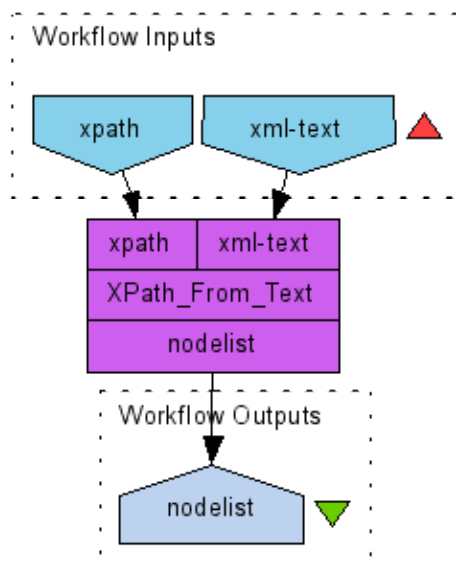
5.4.1 การดึงอิลิเมนต์ข้อมูลในเอกสาร Scufi ของเวิร์คโฟลว์

การดึงอิลิเมนต์ของข้อมูลต่างๆในเอกสาร Scufi ของเวิร์คโฟลว์สามารถใช้ชุดคำสั่งของ XPath เข้าไปดึงอิลิเมนต์ของข้อมูลที่ต้องการได้ ตารางที่ 5.1 แสดงชุดคำสั่ง XPath ในการดึงอิลิเมนต์ข้อมูลที่จะนำมาใช้ในการทำงานตรวจสอบเวิร์คโฟลว์ และโปรแกรมทาวเวอร์น่าก็มีบริการท้องถิ่นชื่อ *'XPath From Text'* (รูปที่ 5.1) รองรับการทำงานนี้ไว้เรียบร้อยแล้ว

ตารางที่ 5.1 XPath expression ที่ใช้ดึงอิลิเมนต์ของข้อมูลในเวิร์คโฟลว์

Extracting	Xpath expression
Processor names	<code>//*[local-name()="Dataset" and (namespace-uri()="http://org.embl.ebi.escience/xScufi-Biomart/0.1.alpha")]/ancestor::s:processor/@name</code>
BioMart processor names	<code>//*[local-name(.)='processor' and (@name=<Processor Name>)] /*[local-name(.)='BioMart']/parent::s:processor/@name</code>
BioMart dataset names	<code>/*[(local-name(.)='processor' and (@name=<A BioMart processor name>)] /*[local-name(.)='BioMart'] /*[local-name(.)='MartQuery'] /*[local-name(.)='Query'] /*[local-name(.)='Dataset']/@name</code>
BioMart dataset filters or attributes	<code>/*[(local-name(.)='processor' and (@name=<A BioMart processor name>)] /*[local-name(.)='BioMart'] /*[local-name(.)='MartQuery'] /*[local-name(.)='Query'] /*[local-name(.)='Dataset'] /*[local-name(.)='Filter' 'Attribute']/@name</code>

รูปที่ 5.7 แสดงการใช้บริการท้องถิ่นอย่างง่าย *'XPath_From_Text'* ซึ่งมีอินพุต 2 อินพุตได้แก่ XPath expression (xpath) ดังแสดงในตารางที่ 5.1 และเอกสาร Scufi ของเวิร์คโฟลว์ (xml-text) โดยจะได้เอาท์พุตเป็นรายการของอิลิเมนต์ (nodelist) ที่ระบุใน XPath expression นั้นๆ



รูปที่ 5.7 การใช้งานบริการ XPath_From_Text

5.4.2 การสืบค้นข้อมูลรายชื่อฟิลเตอร์และแอตทริบิวต์ที่ทันสมัย

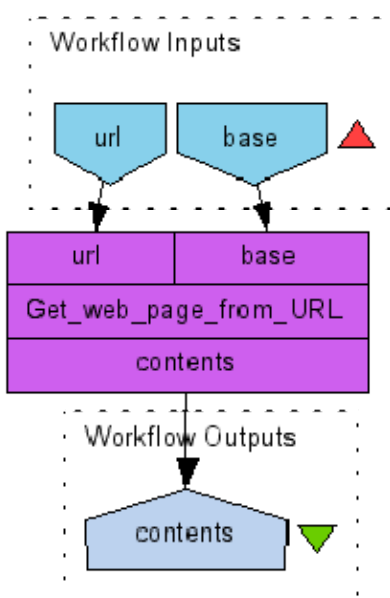
การสืบค้นข้อมูลรายชื่อฟิลเตอร์และแอตทริบิวต์ที่ทันสมัย สามารถทำได้โดยการค้นหาฐานข้อมูลรีจิสทรีกลางของบริการ ไบโอมาร์ท ชุดเมตาดาตาของดาตาเซตที่ให้บริการสามารถเรียกได้จากการเติมพารามิเตอร์ต่อท้าย URL หลัก (Base URL) พารามิเตอร์แรกคือเครื่องหมาย '?' จากนั้นพารามิเตอร์อื่นๆคือเครื่องหมาย '&' เมตาดาตาของดาตาเซตเข้าถึงได้โดยการร้องขอด้วยคำสั่ง *GET* ของ HTTP ดังตารางที่ 5.2

ตารางที่ 5.2 รูปแบบ URL สำหรับการสืบค้นรีจิสทรีของบริการไบโอมาร์ท

ประเภท (Type)	การอิมพลีเมนต์ (Implementation)
Base URL	http://www.biomart.org
ฟิลเตอร์รีจิสทรี (Filter registry)	/biomart/martservice?type=filters&virtualschema=default&dataset =<datasetname>
แอตทริบิวต์รีจิสทรี (Attribute registry)	/biomart/martservice?type=attributes&virtualschema=default&dataset =<datasetname>

ตัวอย่างการค้นหารายชื่อแฉัตริบิวต์ต่างๆของดาตาเซต เช่น ดาตาเซตชื่อ *hsapiens_gene_ensembl* จะมี Base URL คือ <http://www.biomart.org> และมีแฉัตริบิวต์รหัสคือ `/biomart/martservice?type=attributes&virtualschema=default&dataset=hsapiens_gene_ensembl`

โปรแกรมทาเวอร์นามิบริการท้องถิ่นชื่อ *Get_web_page_from_URL* รองรับการทำงานในการร้องขอ *GET* ข้อมูลผ่านเครือข่ายอินเทอร์เน็ตของ HTTP รูปที่ 5.8 แสดงการใช้งานอย่างง่ายของบริการ โดยมีการรับอินพุตสอง อย่างได้แก่ URL หลัก (base) และพารามิเตอร์ (url) จะได้อาท์พุตคือรายชื่อของฟิลด์และแฉัตริบิวต์ต่างๆตามที่ระบุในพารามิเตอร์



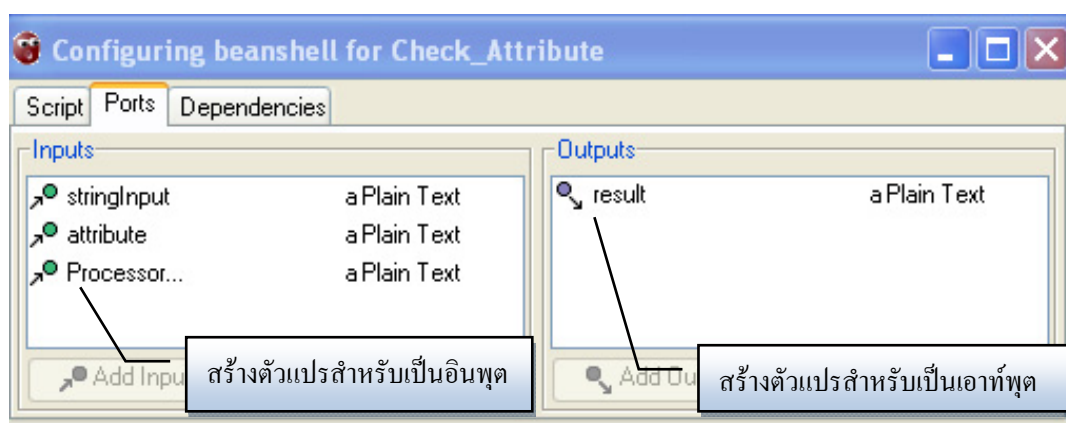
รูปที่ 5.8 ตัวอย่างการใช้งานบริการ *Get_web_page_from_URL*

5.4.3 การเปรียบเทียบข้อมูลในเวิร์คโฟลว์กับบริจิสทรีของไบโอมาร์ท

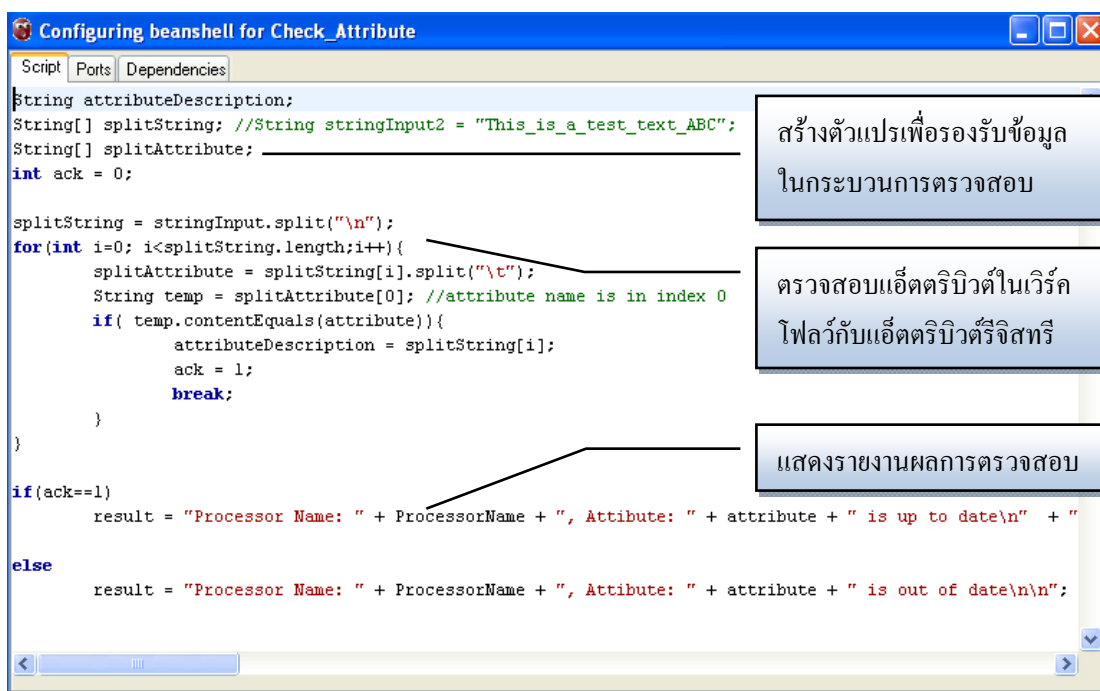
เมื่อสามารถดึงอิลิเมนต์ของข้อมูลต่างๆที่ต้องการจากเอกสาร Scuff ของเวิร์คโฟลว์ และสามารถสืบค้นข้อมูลฟิลด์และแฉัตริบิวต์ต่างๆที่ทันสมัยจากฐานข้อมูลริจิสทรีกลางของบริการไบโอมาร์ทได้แล้ว จึงนำข้อมูลสองส่วนนี้มาเปรียบเทียบกันเพื่อหาว่าฟิลด์และแฉัตริบิวต์ใดที่ได้ล้าสมัยไปแล้วหรือยังคงทันสมัยอยู่ กระบวนการทำงานในส่วนนี้ต้องใช้ภาษาคอมพิวเตอรืในการพัฒนาโปรแกรมเข้ามาช่วยเหลือ โดยภายใต้การพัฒนาเวิร์คโฟลว์ในสิ่งแวดล้อมของโปรแกรมทาเวอร์นามิบริการท้องถิ่นชื่อว่าบิเนลล์ซึ่งเป็นลักษณะเซลล์ทำหน้าที่

ตีความสคริปต์ของภาษา Java ได้ ดังนั้นจึงสามารถเขียนอัลกอริทึมในรูปแบบของสคริปต์ฝังไว้ในบีนเชลล์ได้

กระบวนการที่ 3.3 จากรูปที่ 5.6 แสดงให้เห็นว่า กระบวนการรับข้อมูล 3 ชนิดคือฟิลด์อร์และแอตทริบิวต์ที่ถูกดึงมาจาก Scuff ของเวิร์คโฟลว์, ข้อมูลรหัสทรีของคาตาเซตที่ฟิลด์อร์และแอตทริบิวต์เป็นสมาชิก และชื่อของบริการที่เรียกใช้คาตาเซตนั้นๆ ทำให้สามารถปรับแต่งเพื่อให้บีนเชลล์รับข้อมูลอินพุตได้ 3 พอร์ตและมี 1 พอร์ตสำหรับรายงานผลดังรูปที่ 5.9



รูปที่ 5.9 การกำหนดอินพุตและเอาต์พุตในบีนเชลล์



รูปที่ 5.10 สคริปต์ของบีนเชลล์ในการเปรียบเทียบแอตทริบิวต์ของบริการไบโอมาร์ท

การกำหนดอินพุตและเอาต์พุตในบีนเชลล์ก็คือ แนวคิดการกำหนดตัวแปรเพื่อรองรับข้อมูลที่ไหลเข้าและไหลออกจากบีนเชลล์ เมื่อกำหนดตัวแปรเรียบร้อยแล้วสามารถเขียนสคริปต์ตามอัลกอริทึมในการเปรียบเทียบฟิลด์และแอตทริบิวต์ที่ออกแบบจาก Data flow diagram ได้ ตัวอย่างสคริปต์ของบีนเชลล์ในการตรวจสอบแอตทริบิวต์แสดงดังรูปที่ 5.10

5.4.4 การปรับแต่งกลไกการทำซ้ำของพอร์ตอินพุตในเวิร์กโฟลว์

วัตถุประสงค์ของการทำซ้ำของพอร์ตอินพุตในโปรแกรมทาวเวอร์นาคือ การจัดการกับกลไกการป้อนข้อมูลอินพุตให้กับกระบวนการการทำงาน ในงานด้านการเขียนโปรแกรมภาษาคอมพิวเตอร์ การจัดการข้อมูลอินพุตเหล่านี้ อาจอยู่ในรูปแบบ List หรือรายการข้อมูลที่ประกอบด้วย Item ต่างๆของข้อมูล จากนั้นใช้วิธีการโปรแกรมแบบลูปในการส่งข้อมูลที่ละ Item เข้าไปประมวลในกระบวนการทำงานที่ต้องการนั้น

จากการออกแบบใน Data flow diagram ในรูปที่ 5.5 และ รูปที่ 5.6 นั้น แสดงให้เห็นว่า แต่ละกระบวนการมีการทำงานแบบวนซ้ำหรือลูป ดังนั้นจึงมีกลไกการป้อนข้อมูลอินพุตให้กับกระบวนการนั้นๆ โดยมีอินพุตคือ เอกสาร Scufi ของเวิร์กโฟลว์ (Workflow Scufi), ชื่อบริการไปโอมาร์ททั้งหมด (Processor name) และชื่อคาตาเซตของบริการนั้นๆ (Dataset name) การทำงานก็จะส่งข้อมูลเข้าไปที่ละบริการหรือที่ละคาตาเซต เพราะหนึ่งบริการจะมีเพียงหนึ่งคาตาเซตเท่านั้น

การจัดการอินพุตของโปรแกรมทาวเวอร์นาค้นมี 2 ลักษณะคือ Dot product หรือหนึ่งพอร์ตต่อหนึ่งพอร์ต (One against one) และ Cross product หรือแบบหลายพอร์ตต่อหลายพอร์ต (All against all) ซึ่งค่าโดยปริยายของโปรแกรมทาวเวอร์นาเป็นแบบ Cross product นั้นไม่สอดคล้องกับลักษณะโครงสร้าง Scufi ของบริการไปโอมาร์ทที่จะตรวจสอบในวิทยานิพนธ์นี้

ตารางที่ 5.3 แสดงตัวอย่างของข้อมูล 2 Lists แต่ละ List ประกอบด้วย 3 Items หากจัดการข้อมูลให้มีการทำซ้ำกันแบบ Cross product จะได้ผลลัพธ์ดังตารางที่ 5.4 และหากจัดการข้อมูลให้มีการทำซ้ำกันแบบ Dot product จะได้ผลลัพธ์ดังตารางที่ 5.5

ดังนั้นในวิธานิพนธ์นี้ บริการไปโอมาร์ททั้งหมด (Processor name) และชื่อดาตาเซตของบริการนั้นๆ (Dataset name) จะต้องกำหนดการทำซ้ำกันแบบ Dot product เพราะเป็นความสัมพันธ์แบบหนึ่งต่อหนึ่ง ในขณะที่เอกสาร Scufi ของเวิร์กโฟลว์ (Workflow Scufi) จะต้องกำหนดการทำซ้ำกันแบบ Cross product เพราะบริการไปโอมาร์ทและดาตาเซตเป็นสมาชิกในเอกสาร Scufi ดังนั้นจึงสามารถปรับแต่งกลไกการทำซ้ำได้ดังรูปที่ 5.11

ตารางที่ 5.3 ตัวอย่างของข้อมูลอินพุตแบบหลายรายการ

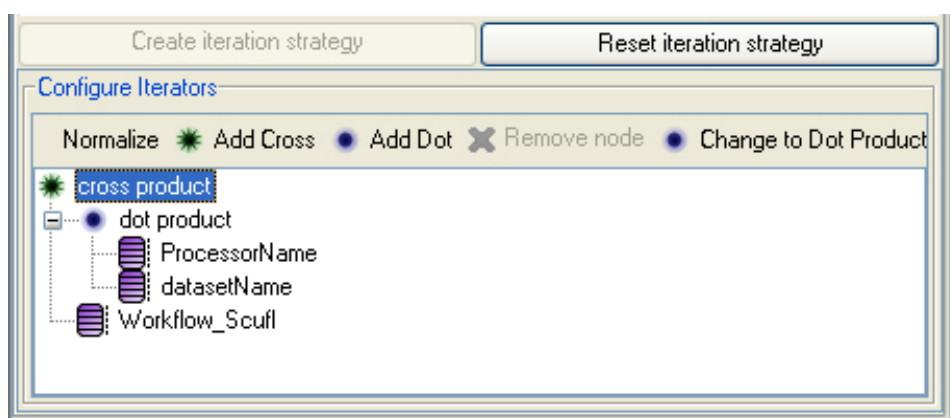
String1	String2
A	B
C	D
E	F

ตารางที่ 5.4 ผลลัพธ์การทำซ้ำแบบ Cross product

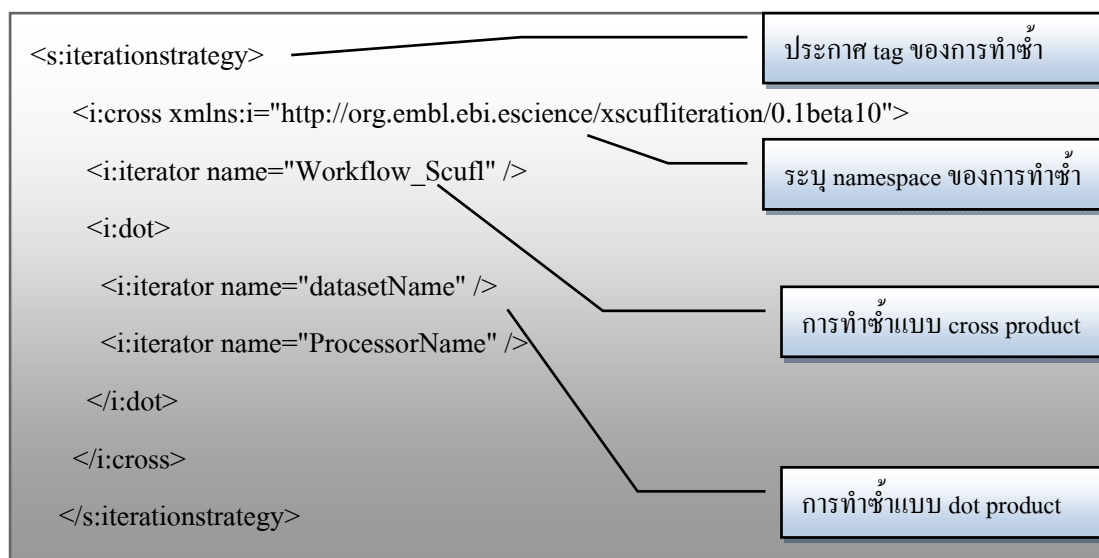
Cross product output
AB
AD
AF
CB
CD
CF
EB
ED
EF

ตารางที่ 5.5 ผลลัพธ์การทำซ้ำแบบ Dot product

Dot product output
AB
CD
EF



รูปที่ 5.11 การปรับแต่งกลไกการทำซ้ำของพอร์ตอินพุตในเวิร์คโฟลว์



รูปที่ 5.12 โครงสร้างของ Scuf1 ในการปรับแต่งการทำซ้ำของพอร์ตอินพุต

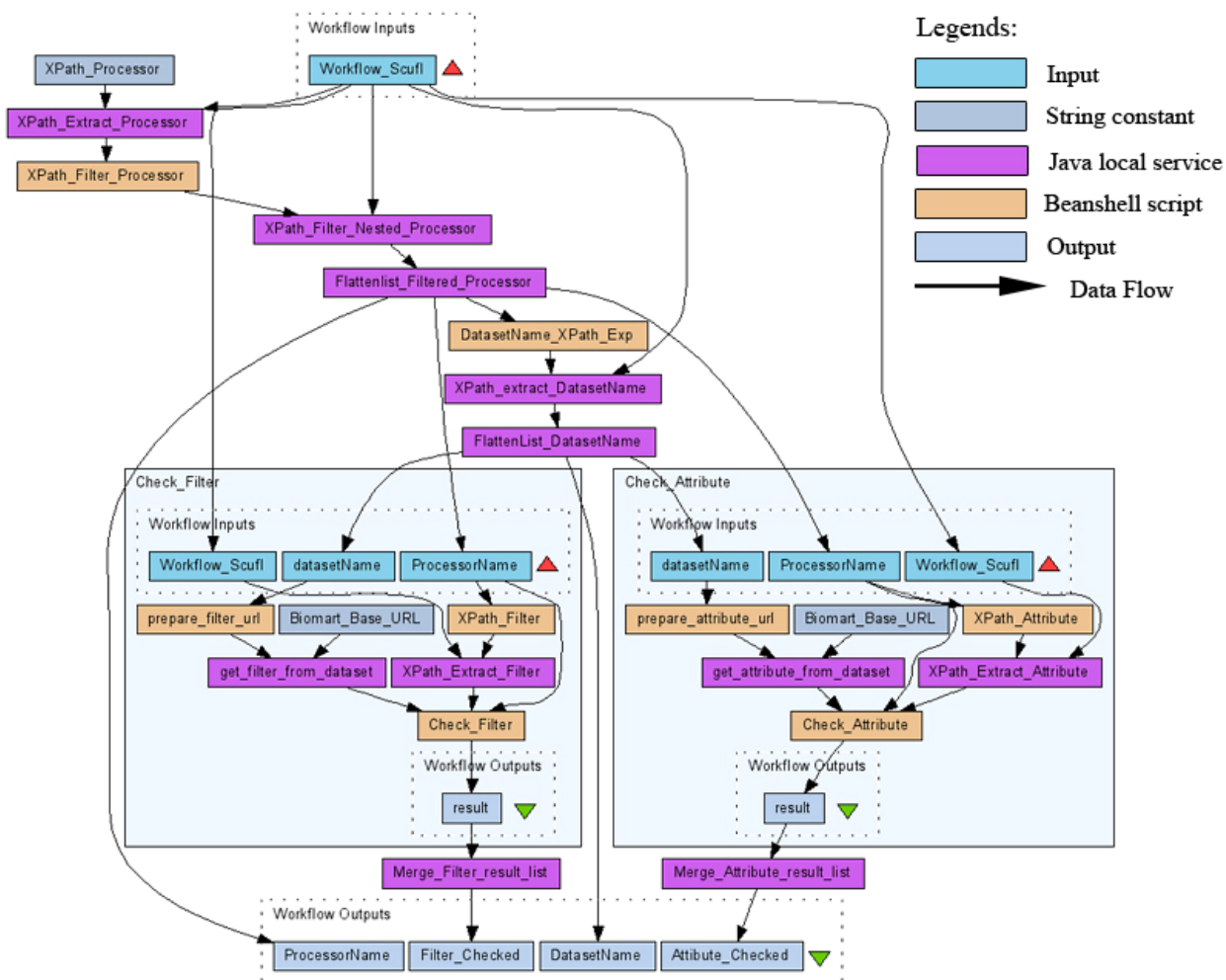
จากการปรับแต่งในรูปที่ 5.11 แสดงความสัมพันธ์ของแต่อินพุตในเวิร์คโฟลว์ เอกสาร Scuf1 ของเวิร์คโฟลว์ (Workflow Scuf1) มีชนิดการทำซ้ำของข้อมูลแบบ Cross product

ส่วนชื่อบริการไปโอมาร์ทั้งหมด (Processor name) และชื่อคาตาเซตของบริการนั้นๆ (Dataset name) มีชนิดการซ้ากันของข้อมูลแบบ Dot product ผลที่ได้คือรูปที่ 5.12 แสดงโครงสร้างของ ScufI ในการปรับแต่งการทำซ้ำของพอร์ตอินพุตเหล่านี้ในเวิร์คโพลว์ ดังนั้นผลลัพธ์ความสัมพันธ์ของกลไกการทำซ้ำโดยใช้ข้อมูลอินพุตทั้ง 3 ดังแสดงในตารางที่ 5.6 สามารถอธิบายได้ว่า เอกสาร ScufI ของเวิร์คโพลว์ใดๆจะประกอบด้วยจำนวนบริการได้หลายบริการ และแต่ละบริการประกอบด้วยคาตาเซตหนึ่งคาตาเซต

ตารางที่ 5.6 ผลลัพธ์กลไกการทำซ้ำของข้อมูลอินพุตในเวิร์คโพลว์

Cross product	Dot product	Dot product
Workflow ScufI	1 st Dataset name	1 st Processor name
Workflow ScufI	2 nd Dataset name	2 nd Processor name

จากการออกแบบและการพัฒนาตามลำดับขั้นตอนดังกล่าวแล้ว สามารถสร้างเวิร์คโพลว์สำหรับการตรวจสอบความถูกต้องของบริการไปโอมาร์ที่ได้ดังรูปที่ 5.13 และกรณีศึกษาของเวิร์คโพลว์นี้ก็คือ เวิร์คโพลว์สำหรับการวิเคราะห์สนธิปของมนุษย์ที่กล่าวถึงการออกแบบและการพัฒนาแล้วในบทที่ 4

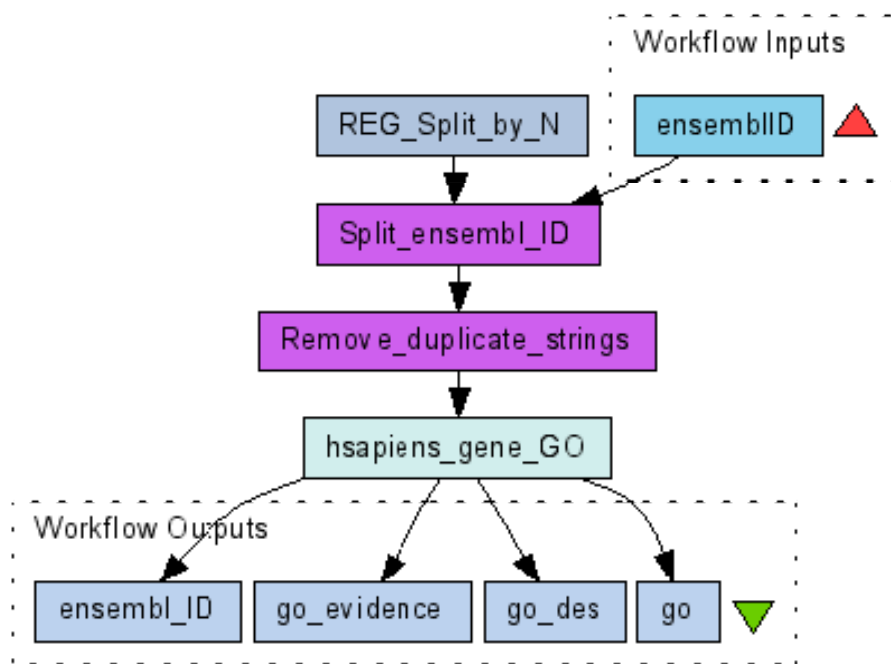


รูปที่ 5.13 เวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้องของบริการไบโอมาร์ท

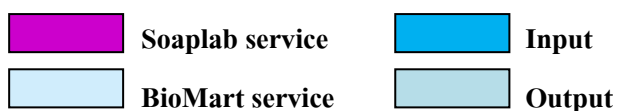
5.5 การตรวจสอบเวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท

ในเบื้องต้นได้สร้างเวิร์คโฟลว์อย่างง่าย ประกอบด้วยบริการไบโอมาร์ทที่ทำหน้าที่ในการคิวรีข้อมูลจำนวน 1 บริการชื่อ *hsapiens_gene_ensembl* แต่เปลี่ยนชื่อเป็น *hsapiens_gene_GO* เพื่อบอกจุดประสงค์การทำงานค้นหา Gene ontology เวิร์คโฟลว์ประกอบด้วย 1 ฟิลเตอร์สำหรับเป็นเงื่อนไขการคิวรีข้อมูลคือ Ensembl ID และ 3 แอ็ตทริบิวต์สำหรับเป็นข้อมูลเอาท์พุทเกี่ยวกับ Gene Ontology คือ Gene ontology, Gene ontology description และ Gene ontology evidence และอีก 1 แอ็ตทริบิวต์คือ Ensembl ID เอง กล่าวคือ Ensembl ID ที่เป็นอินพุทของฟิเตอร์คิวรีข้อมูลในคาตาเซ็ต *hsapiens_gene_ensembl* โดยจะคิวรี Ensembl ID มาด้วยเช่นกัน

เพื่อใช้ตรวจสอบความถูกต้องของบริการ เวิร์คโฟลว์ทดสอบอย่างง่ายดังรูปที่ 5.14 และตัวอย่างการทำงาน
 ของบริการท้องถิ่น Java ดังตารางที่ 5.7 โดยมีโครงสร้างภาษา Scufi ของเวิร์คโฟลว์ดังรูปที่
 5.15 และมีผลการตรวจสอบด้วยเวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ทดังนี้



Legends:



รูปที่ 5.14 เวิร์คโฟลว์อย่างง่ายสำหรับการทดลอง

ตารางที่ 5.7 ตัวอย่างการทำงานของบริการท้องถิ่น Java

ชื่อบริการ	หน้าที่	ตัวอย่างอินพุต	เอาต์พุต
Spilt_ensembl_ID	แยก string	“ENSG00000117242 ENSG00000158828 ENSG00000158828”	“ENSG00000117242” “ENSG00000158828” “ENSG00000158828”
Remove_duplicate_strings	ลบ string ที่ซ้ำซ้อน	“ENSG00000117242” “ENSG00000158828” “ENSG00000158828”	“ENSG00000117242” “ENSG00000158828”

จากรายงานในรูปแบบที่ 5.16 ทำให้รู้ว่า แอ็ตทริบิวต์ทุกตัวของการทำงานในส่วน Gene ontology นั้นได้ล้าสมัยไปแล้ว (Out-of-date) ซึ่งสามารถตรวจสอบโดยวิธีการแบบตรงไปตรงมาว่า แอ็ตทริบิวต์เหล่านี้ได้ล้าสมัยจริงหรือไม่ โดยตรวจสอบที่หน้าตาการปรับแต่งของบริการไบโอ-มาร์ท (Configuration screen) แล้วค้นหาที่หน้าจอนั้นว่า แอ็ตทริบิวต์ดังกล่าวปรากฏบนหน้าตาหรือไม่ ซึ่งแอ็ตทริบิวต์ที่ล้าสมัยไปแล้ว จะไม่ปรากฏที่หน้าจอการปรับแต่งของบริการไบโอ-มาร์ทอีก เนื่องจากการเรียกหน้าตาการปรับแต่งของบริการไบโอ-มาร์ท ณ ขณะเวลาใดๆที่กำลังใช้งานโปรแกรมทาวเวอร์นา โปรแกรมจะใช้ข้อมูลของบริการที่ทันสมัยที่สุด ณ เวลานั้นๆ

นอกจากนี้เวิร์คโฟลว์ที่พัฒนา ยังสามารถให้ข้อมูลที่จำเป็นของฟิลเตอร์และแอ็ตทริบิวต์ที่ทันสมัยสำหรับการตรวจสอบในอนาคตได้อีกด้วย โดยให้อาท์พุตเป็นรายงานดังรูปที่ 5.17

จากรูปที่ 5.17 รายงานข้างต้นทำให้ทราบว่า ฟิลเตอร์ชื่อ *ensembl_gene_id* มีคำอธิบายว่าเป็นเงื่อนไขการค้นหาข้อมูลโดยสามารถรองรับเป็นรายชื่อของ Ensembl Gene ID ได้ ส่วนแอ็ตทริบิวต์ชื่อ *ensembl_gene_id* ก็คือ Ensembl Stable ID ของยีนนั้นๆ ซึ่งจะปรากฏอยู่ในหน้า *Feature* ของหน้าจอการปรับแต่งของบริการไบโอ-มาร์ท ซึ่งสามารถตรวจสอบความถูกต้องของรายงานการทำงานของเวิร์คโฟลว์โดยวิธีการแบบ Manual โดยตรงได้

จากการตรวจสอบโดยวิธีการแบบ Manual แอ็ตทริบิวต์ที่ล้าสมัยจะไม่ปรากฏอยู่ในหน้าจอการปรับแต่งแล้วอันเป็นผลมาจากการปรับปรุงบริการไบโอ-มาร์ท จึงสามารถสรุปได้ว่าบริการไบโอ-มาร์ทที่ทดสอบนี้ ได้ล้าสมัยไปแล้ว



รูปที่ 5.15 โครงสร้างภาษา Scufl ของเวิร์กโฟลว์ตัวอย่างง่ายสำหรับการทดลอง

Processor Name: hsapiens_gene_GO, Attribute: go is out-of-date
Processor Name: hsapiens_gene_GO, Attribute: evidence_code is out-of-date
Processor Name: hsapiens_gene_GO, Attribute: go_description is out-of-date

รูปที่ 5.16 รายงานการตรวจสอบบริการไบโอมาร์ทในส่วนของฟิลด์ข้อมูลที่ล้าสมัย

Processor Name: hsapiens_gene_GO, Filter: ensembl_gene_id is up-to-date
Detail: Filter to include genes with supplied list of Ensembl Gene IDs
Processor Name: hsapiens_gene_GO, Attribute: ensembl_gene_id is up-to-date
Detail: Ensembl Stable ID of the Gene, feature_page

รูปที่ 5.17 รายงานการตรวจสอบบริการไบโอมาร์ทในส่วนของฟิลด์ข้อมูลที่ยังทันสมัย

5.6 สรุป

บทนี้กล่าวถึงแนวคิดในการออกแบบ การพัฒนา และการทดสอบเวิร์ค โฟลว์ สำหรับตรวจสอบการเปลี่ยนแปลงของบริการไบโอมาร์ท ในเวิร์ค โฟลว์สำหรับวิเคราะห์ของมนุษย์แบบอัตโนมัติ และเวิร์ค โฟลว์ที่พัฒนาขึ้นใช้บริการท้องถิ่นที่โปรแกรมทาเวอร์น่าเตรียมไว้ให้ จาก การทดสอบด้วยเวิร์ค โฟลว์กรณีศึกษาอย่างง่ายพบว่า เวิร์ค โฟลว์สำหรับตรวจสอบการเปลี่ยนแปลงของบริการไบโอมาร์ทซึ่งสามารถระบุได้ว่า บริการและฟิลด์หรือแอตทริบิวต์ตัวใดที่ล้าสมัยไปแล้ว เพราะนั่นหมายความว่าบริการไบโอมาร์ทนั้นๆ ได้ล้าสมัยไปแล้ว ในบทต่อไปกล่าวถึงการทดสอบและผลการทดสอบของเวิร์ค โฟลว์ สำหรับตรวจสอบความถูกต้องของบริการไบโอมาร์ทที่พัฒนาขึ้นในบทที่ 4 เป็นกรณีศึกษา

บทที่ 6

การตรวจสอบบริการไบโอมาร์ทในเวิร์คโฟลว์สำหรับวิเคราะห์สลับของมนุษย์

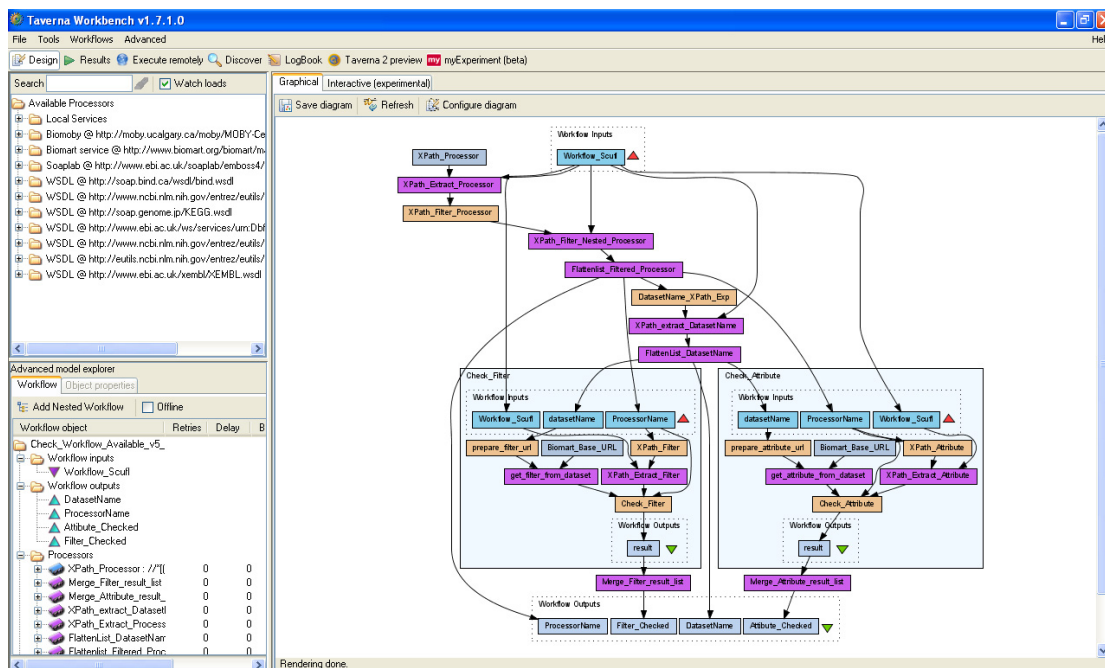
บทนำ

จากบทที่ 4 ได้ทำการทดสอบการทำงานของเวิร์คโฟลว์สำหรับวิเคราะห์สลับของมนุษย์ซึ่งใช้ SNP IDs เป็นอินพุตให้กับเวิร์คโฟลว์จำนวน 1,000 IDs ปรากฏว่าเวิร์คโฟลว์ใช้เวลาทำงานประมาณ 14 ชั่วโมง และผลลัพธ์ที่ได้ก็ไม่ถูกต้องตรงตามความต้องการ ทำให้ต้องเสียเวลาเป็นชั่วโมงๆหรือเป็นวัน ส่วนในบทที่ 5 ได้นำเสนอวิธีการออกแบบและพัฒนาเวิร์คโฟลว์สำหรับตรวจสอบความถูกต้องของเวิร์คโฟลว์ก่อนการทำงานจริง โดยจะมุ่งเน้นไปที่บริการไบโอมาร์ทเพื่อลดเวลาในการดีบั๊กหรือตรวจสอบเวิร์คโฟลว์ได้อย่างมีนัยสำคัญ ซึ่งผลการทดสอบเวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท กับเวิร์คโฟลว์กรณีศึกษาอย่างง่าย พบว่าเวิร์คโฟลว์ที่พัฒนาสามารถระบุบริการไบโอมาร์ทที่ล่าสมัยได้ถูกต้อง ส่วนในบทนี้จะกล่าวถึงการตรวจสอบบริการไบโอมาร์ท [60] โดยใช้เวิร์คโฟลว์สำหรับวิเคราะห์สลับของมนุษย์ซึ่งมีความซับซ้อนกว่าหลายเท่าเป็นกรณีศึกษา จะกล่าวถึงกระบวนการหลังการตรวจสอบ ข้อเสนอแนะและสรุป

6.1 การตรวจสอบเวิร์คโฟลว์สำหรับวิเคราะห์สลับของมนุษย์

จากการทดสอบการทำงานของเวิร์คโฟลว์ เพื่อการตรวจสอบบริการไบโอมาร์ทกับเวิร์คโฟลว์กรณีศึกษาอย่างง่ายในการค้นหา Gene ontology ในบทที่ 5 แล้วนั้น พบว่าเวิร์คโฟลว์ที่พัฒนาขึ้นมาสามารถทำงานได้ถูกต้อง โดยสามารถระบุฟิลเตอร์และแอตทริบิวต์ที่ล่าสมัยและให้ข้อมูลที่สำคัญของฟิลเตอร์และแอตทริบิวต์ที่ทันสมัยได้ ดังนั้นจึงตรวจสอบเวิร์คโฟลว์สำหรับการวิเคราะห์สลับของมนุษย์ ซึ่งเป็นเวิร์คโฟลว์ที่ใช้มีบริการและการเชื่อมโยงข้อมูลที่ซับซ้อนและใช้ฐานข้อมูลที่หลากหลายดังที่กล่าวแล้วในบทที่ 4 และสามารถเข้าถึงได้ตามที่อยู่ URL ดังนี้ <http://www.myexperiment.org/workflows/610> ขั้นตอนการตรวจสอบมีดังนี้

- 1) เปิดเวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้องของบริการไบโอมาร์ทด้วยโปรแกรมทาเวอร์น่าดังรูปที่ 6.1



รูปที่ 6.1 เวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท

2) เริ่มการตรวจสอบบริการไบโอมาร์ท ในเวิร์คโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์โดยการสั่ง Run เวิร์คโฟลว์แล้วจะได้หน้าต่างดังรูปที่ 6.2 จะเห็นว่าหน้าต่างด้านบนแสดงเวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท และหน้าต่างด้านล่างเป็นพื้นที่สำหรับใส่อินพุตให้เวิร์คโฟลว์ทำงาน อินพุตในวิทยานิพนธ์นี้ก็คือ เอกสาร Scum ของเวิร์คโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์

3) เมื่อใส่อินพุตพร้อมแล้วก็สามารถสั่ง Run เวิร์คโฟลว์ได้ รูปที่ 6.3 แสดงเวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้องบริการไบโอมาร์ทกำลังทำงาน บล็อกที่มีสีเขียวแสดงสถานะการทำงานว่า บริการนี้ทำงานสำเร็จโดยไม่มีข้อผิดพลาด บล็อกที่มีสีม่วงแสดงสถานะการทำงานว่า บริการนี้กำลังทำงาน และบล็อกที่มีสีเหลืองบอกให้รู้ว่าบริการนี้มีการทำงานวนซ้ำหรือลูป จากรูปที่ 6.3 จะเห็นว่าคือบล็อกของบริการตรวจสอบฟิลเตอร์และแอ็คตริบิวต์

4) เมื่อเวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้อง ของบริการไบโอมาร์ททำงานเสร็จสิ้น จะได้ผลลัพธ์เป็นรายงานการตรวจสอบฟิลเตอร์และแอ็คตริบิวต์ รูปที่ 6.4 แสดงหน้าต่างรายงานผลการตรวจสอบฟิลเตอร์ของบริการไบโอมาร์ท โดยรายงานดังแสดงในรูปที่ 6.5 ซึ่งจะเห็นว่าฟิลเตอร์ทุกฟิลเตอร์ของบริการไบโอมาร์ทในเอกสาร Scum อินพุตที่ใช้ในการตรวจสอบมีความทันสมัยและให้ข้อมูลรายละเอียดที่สำคัญของฟิลเตอร์นั้นๆ

5) รูปที่ 6.6 แสดงหน้าต่างรายงานผลการตรวจสอบแอ็ดตริบิวต์ของบริการไปโอ-มาร์ท โดยรายงานแสดงดังรูปที่ 6.7 จะเห็นว่าแอ็ดตริบิวต์ทุกตัวของบริการไปโอมาร์ทล้าสมัยเสียแล้วและทำให้ไม่มีข้อมูลใดๆที่เกี่ยวข้องกับแอ็ดตริบิวต์ที่ล้าสมัยนั้นๆอยู่ที่บริการไปโอมาร์ทอีก

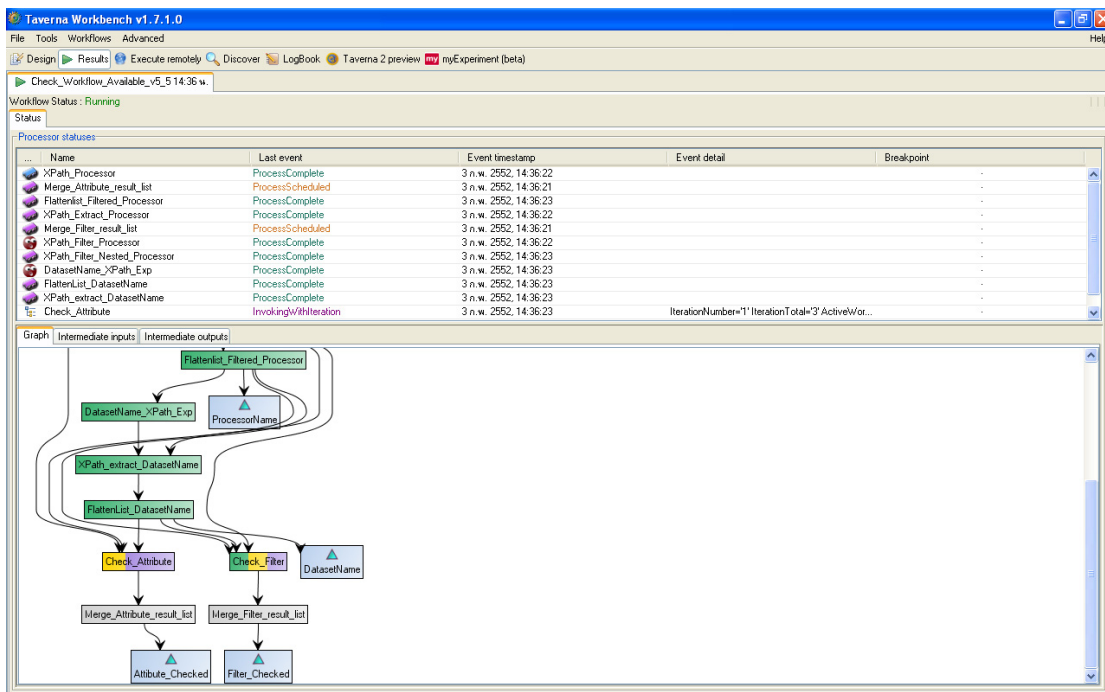
The screenshot displays a workflow execution window titled "Run workflow: Check_Workflow_Available_v5_5". The main area shows a complex flowchart with various nodes and arrows representing the workflow process. Below the flowchart, there is a section labeled "Inputs" with a toolbar and a list of input documents. The "Workflow_Scuff" document is selected, showing its XML content. Two callout boxes provide additional context: one points to the flowchart and the other points to the XML code.

เวิร์กโฟลว์สำหรับการตรวจสอบบริการไปโอมาร์ท

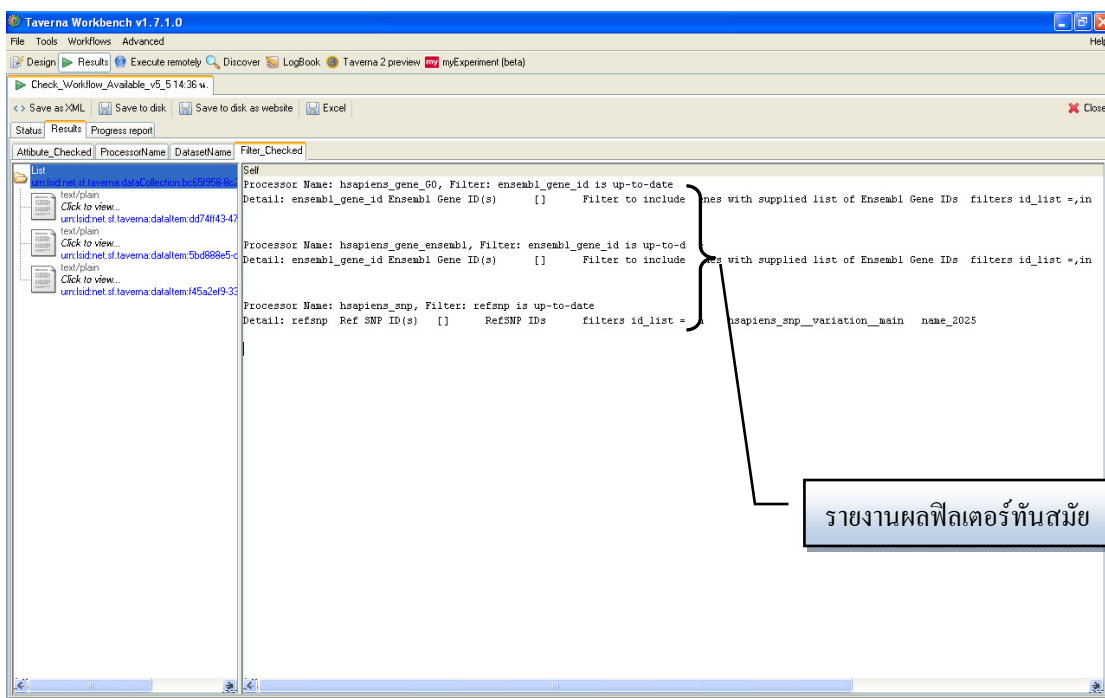
อินพุตคือเอกสาร Scuff ของเวิร์กโฟลว์ในงานวิจัยเภสัชพันธุศาสตร์

```
<?xml version="1.0" encoding="UTF-8"?>
<s:scuff xmlns:s="http://org.enbi.esience/xscuff/0.1.alpha" version="0.2" log="0">
<s:workflowdescription lsd="urn:lsid:net.sf.taverna.wd.definition:16173a30-74a0-41ce-add5b7471c21f3d" author="" title="BioMart_hsapiens_gene_ensembl_variation" />
<s:processor name="REG_Split_by_C" boring="true">
<s:stringconstant </s:stringconstant>
</s:processor>
<s:processor name="REG_Merge_by_C" boring="true">
<s:stringconstant </s:stringconstant>
</s:processor>
<s:processor name="Flatten_list">
<s:local>
org.enbi.esience.scuffworkers.java.FlattenList
<s:extensions>
<s:flattenlist s:depth="2" />
</s:extensions>
</s:local>
</s:processor>
<s:processor name="Flatten_Ensembl_report">
<s:local>
org.enbi.esience.scuffworkers.java.FlattenList
<s:extensions>
<s:flattenlist s:depth="2" />
</s:extensions>
</s:local>
</s:processor>
</s:scuff>
```

รูปที่ 6.2 หน้าต่างโปรแกรมพร้อมที่จะตรวจสอบเวิร์กโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์



รูปที่ 6.3 เวิร์กโฟลว์กำลังตรวจสอบบริการไบโอมาร์ท



รูปที่ 6.4 ผลการตรวจสอบฟิลเตอร์ของบริการไบโอมาร์ทในโปรแกรมทาวเวอร์น่า

Processor Name: *hsapiens_gene_GO*, Filter: *ensembl_gene_id* is up-to-date
 Detail: *ensembl_gene_id* Ensembl Gene ID(s) [] Filter to include genes with supplied list of Ensembl Gene IDs filters *id_list* =,in
hsapiens_gene_ensembl_gene_main *stable_id_1023*

Processor Name: *hsapiens_gene_ensembl*, Filter: *ensembl_gene_id* is up-to-date
 Detail: *ensembl_gene_id* Ensembl Gene ID(s) [] Filter to include genes with supplied list of Ensembl Gene IDs filters *id_list* =,in
hsapiens_gene_ensembl_gene_main *stable_id_1023*

Processor Name: *hsapiens_snp*, Filter: *refsnp* is up-to-date
 Detail: *refsnp* Ref SNP ID(s) [] RefSNP IDs filters *id_list* =,in
hsapiens_snp_variation_main *name_1023*

รูปที่ 6.5 รายงานการตรวจสอบฟิลเตอร์ของบริการไบโอมาร์ท

The screenshot shows the Taverna Workbench interface with a workflow execution report. The report lists several processors and their attributes. Two boxes on the right highlight specific results:

- รายงานผลเอ็ดตรีบีวต์ ล้าสมัย** (Outdated EDR Report): This box points to three processors with attributes: *go* is out-of-date, *evidence_code* is out-of-date, and *go_description* is out-of-date.
- รายงานเอ็ดตรีบีวต์ ทันสมัย** (Up-to-date EDR Report): This box points to several processors with attributes: *ensembl_gene_id* is up-to-date, *ensembl_gene_id* is up-to-date, *description* is up-to-date, *chromosome_name* is up-to-date, *start_position* is up-to-date, *end_position* is up-to-date, and *strand* is up-to-date.

รูปที่ 6.6 ผลการตรวจสอบเอ็ดตรีบีวต์ของบริการไบโอมาร์ทในโปรแกรมทาวเวอร์น่า

Processor Name: *hsapiens_gene_GO*, Attribute: *go* is out-of-date
 Processor Name: *hsapiens_gene_GO*, Attribute: *evidence_code* is out-of-date
 Processor Name: *hsapiens_gene_GO*, Attribute: *go_description* is out-of-date

รูปที่ 6.7 รายงานการตรวจสอบเอ็ดตรีบีวต์ล้าสมัยของบริการไบโอมาร์ท

ตารางที่ 6.1 แสดงคุณสมบัติและผลการทดสอบเวิร์คโฟลว์สำหรับวิเคราะห์ SNP ของมนุษย์ การทดสอบวัดผลการทดสอบ 10 ครั้ง เวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้องของบริการไบโอมาร์ก ใช้เวลาในการทำงานเฉลี่ย 26.5 วินาทีในการตรวจสอบบริการฟิลเตอร์ และแอ็คตริวิตีได้อย่างถูกต้อง ซึ่งหากไม่ตรวจสอบก่อนจะใช้เวลาจนถึง 14 ชั่วโมงกว่าจะรู้ว่าเวิร์คโฟลว์ดังกล่าวทำงานผิดพลาด

จากรูปที่ 6.4 จะเห็นว่าฟิลเตอร์ทุกตัวมีความทันสมัย โดยรายงานการตรวจสอบฟิลเตอร์ของบริการไบโอมาร์กในรูปที่ 6.5 ซึ่งอธิบายข้อมูลตามที่ได้ไปดึงจากรีจิสทรีของบริการนั้นๆ ซึ่งข้อมูลเหล่านี้จะเป็นประโยชน์ในการอ้างอิงสำหรับการตรวจสอบในอนาคตต่อไป เพราะในอนาคตฟิลเตอร์เหล่านี้อาจจะถูกปรับปรุงได้เช่นกัน

จากรูปที่ 6.6 จะเห็นว่าแอ็คตริวิตีบางตัวที่ปรากฏในเอกสาร Scufi อินพุตได้ล้าสมัยไปแล้ว โดยในรายงานการตรวจสอบแสดงดังรูปที่ 6.7 พบว่าบริการที่ล้าสมัยไปแล้วก็คือบริการไบโอมาร์กที่ทำหน้าที่ในการค้นหา Gene ontology ซึ่งเวิร์คโฟลว์ที่พัฒนาขึ้นนี้สามารถระบุได้ถูกต้องตรงตามสมมุติฐานและยังแสดงให้เห็นว่าเวิร์คโฟลว์สามารถตรวจสอบกับ Nested Workflow ได้อีกด้วย

ตารางที่ 6.1 คุณสมบัติและผลการทดสอบเวิร์คโฟลว์เภสัชพันธุศาสตร์

รายการ	จำนวน	หมายเหตุ
บริการทั้งหมด	55	-
ดาตาลิงค์ทั้งหมด	161	-
บริการที่เกี่ยวข้องกับบริการไบโอมาร์ท	17	เป็นบริการไบโอมาร์ท 3 บริการ
ดาตาลิงค์ที่เกี่ยวข้องกับบริการไบโอมาร์ท	70	-
ฟิลเตอร์	3	เป็นฟิลเตอร์ของบริการไบโอมาร์ท 3 บริการ
แอ็ตทริบิวต์	22	เป็นแอ็ตทริบิวต์ของบริการไบโอมาร์ท 3 บริการ
ข้อมูลอินพุต	1000	SNP IDs
เวลาทำงาน	14.00.50	ชั่วโมง
เวลาในการตรวจสอบ	26.5	วินาที / ทุกบริการทันสมัยยกเว้น Gene ontology

6.2 กระบวนการหลังการตรวจสอบ

เมื่อได้ตรวจสอบความถูกต้องของเวิร์คโฟลว์ และพบว่ามึบริการที่ล้ำสมัยจึงจำเป็นต้องปรับปรุงแก้ไขให้เวิร์คโฟลว์มีความถูกต้องสมบูรณ์ พร้อมทั้งจะไปใช้งานที่เกี่ยวข้องกับงานวิจัยได้ตามความต้องการ ดังนั้นกระบวนการหลังการตรวจสอบจึงเป็นเรื่องที่จำเป็นเพื่อให้เวิร์คโฟลว์สามารถทำงานได้ถูกต้อง

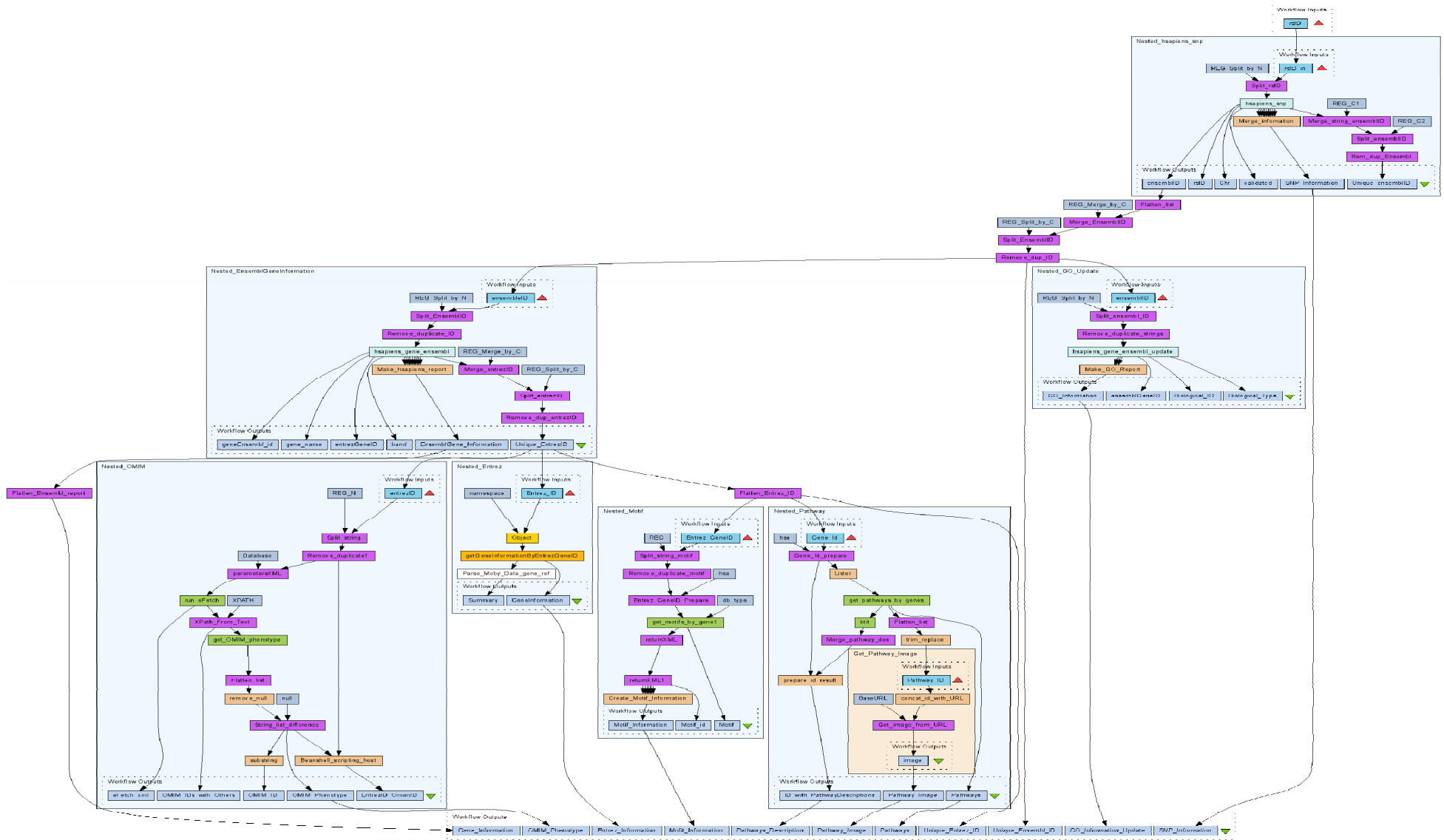
6.2.1 การปรับปรุงและการทดสอบซ้ำเวิร์คโฟลว์สำหรับวิเคราะห์สลิปของมนุษย์

จากผลการตรวจสอบ ทราบว่าแอ็ตทริบิวต์บางตัวของบริการไบโอมาร์ทชื่อ *hsapiens_gene_GO* ซึ่งมีการเรียกใช้ข้อมูลจากดาตาเซ็ต *hsapiens_gene_ensembl* ในการค้นหา

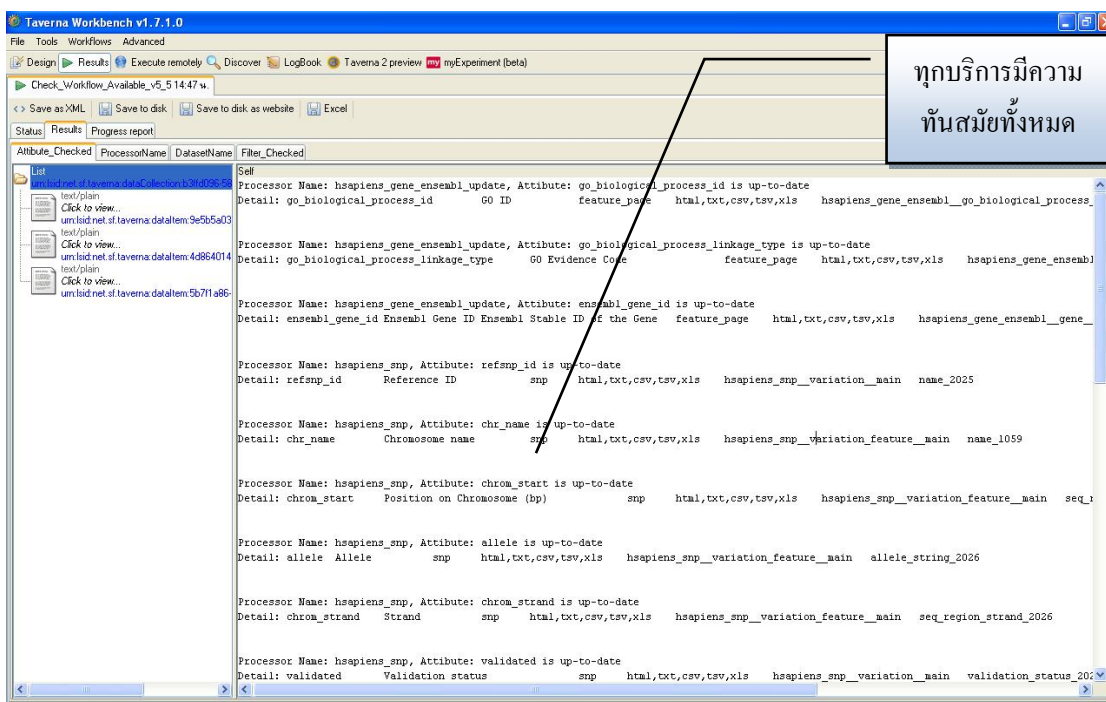
ข้อมูล Gene ontology เป็นบริการที่ล้าสมัยแล้ว สำหรับวิธีการปรับปรุงก็คือ จะต้องเอาบริการนี้ออกจากไดอะแกรมของเวิร์คโฟลว์แล้วทำการเพิ่มอันใหม่เข้ามา โดยอาจจะยังเป็นชื่อเดิมอยู่คือ *hsapiens_gene_ensembl* เพราะรายชื่อของบริการที่นำเสนอต่อผู้ใช้ใน ณ ขณะเวลาใดๆ ย่อมถือว่าทันสมัยที่สุดในขณะนั้น เมื่อเพิ่มบริการใหม่เข้ามาแล้วก็ปรับแต่งเลือกข้อมูลในส่วนที่ขาดหายไป ในกรณีนี้คือ Gene ontology และเชื่อมต่อการเชื่อมโยงข้อมูลกับพอร์ตอื่นๆต่อไป ในกรณีที่บริการที่เปลี่ยนใหม่แล้วนี้ ไม่มีฟิลเตอร์หรือแอ็คตริบิวต์ที่ต้องการ ก็จำเป็นจะต้องค้นหาบริการอื่นๆที่ให้ผลลัพธ์ได้ตามความต้องการในหมวดหมู่ที่เกี่ยวข้อง

สำหรับเหตุผลที่ต้องเอาบริการที่ล้าสมัยออกและเพิ่มบริการใหม่มาแทนที่ ทั้งๆที่บริการเหล่านี้มีชื่อและคาดเดาซึ่งกัน เนื่องจากเราเข้าไปปรับแต่งบริการในหน้าต่างการปรับแต่งของบริการไบโอมาร์กนั้นจะไม่ส่งผลกระทบต่อใดๆ กับฟิลเตอร์และแอ็คตริบิวต์ที่ล้าสมัยเลย เพราะว่าฟิลเตอร์และแอ็คตริบิวต์เหล่านั้นได้ถูกปรับปรุงไปแล้ว โดยอาจจะถูกผู้ให้บริการเอาออกจากบริการไปเลยหรือเปลี่ยนไปเป็นชื่อใหม่ก็ได้แล้วแต่กรณี ดังนั้นในหน้าต่างการปรับแต่งจึงไม่สามารถยกเลิกการเลือกฟิลด์ข้อมูลเหล่านี้ได้ ทำให้ไม่มีผลกระทบต่อเอกสาร Scufi ของเวิร์คโฟลว์แต่อย่างใด รูปที่ 6.8 แสดงเวิร์คโฟลว์สำหรับวิเคราะห์สลับของมนุษย์ที่ได้รับการแก้ไขเรียบร้อยแล้ว และสามารถเข้าถึงได้ทางเว็บไซต์ <http://www.myexperiment.org/workflows/612>

เมื่อแก้ไขเวิร์คโฟลว์เรียบร้อยแล้ว จากนั้นทำการตรวจสอบซ้ำอีกครั้งหนึ่ง ผลการทดสอบคือบริการไบโอมาร์กทุกบริการมีความทันสมัยทั้งหมดแล้ว และพร้อมทำงานได้ตามวัตถุประสงค์ ผลการตรวจสอบในโปรแกรมทาวเวอร์นาแสดงดังรูปที่ 6.9 และตัวอย่างรายงานผลการตรวจสอบของแอ็คตริบิวต์ซึ่งทันสมัยทุกแอ็คตริบิวต์และให้ข้อมูลที่เกี่ยวข้องได้ถูกต้อง แสดงดังรูปที่ 6.10



รูปที่ 6.8 เวิร์กโฟลว์สำหรับวิเคราะห์สปีชของมนุษย์ที่ได้รับการแก้ไขแล้ว



รูปที่ 6.9 การตรวจสอบเวิร์กโฟลว์สำหรับวิเคราะห์ SNP ของมนุษย์หลังการปรับแก้แล้ว

Processor Name: hsapiens_gene_ensembl_update, Attribute: go_biological_process_id is up-to-date
Detail: go_biological_process_id GO ID feature_page html,txt, csv, tsv, xls hsapiens_gene_ensembl_go_biological_process_dm dbprimary_acc_1074

Processor Name: hsapiens_gene_ensembl_update, Attribute: go_biological_process_linkage_type is up-to-date
Detail: go_biological_process_linkage_type GO Evidence Code feature_page html,txt, csv, tsv, xls hsapiens_gene_ensembl_go_biological_process_dm linkage_type_1024

Processor Name: hsapiens_gene_ensembl_update, Attribute: ensembl_gene_id is up-to-date
Detail: ensembl_gene_id Ensembl Gene ID Ensembl Stable ID of the Gene feature_page html,txt, csv, tsv, xls hsapiens_gene_ensembl_gene_main stable id 1023

รูปที่ 6.10 ตัวอย่างรายงานผลการตรวจสอบของแอดตริบิวต์ซึ่งมีความทันสมัยทุกตัว

6.2.2 การตรวจสอบความถูกต้องของผลลัพธ์

ผลลัพธ์ของเวิร์คโฟลว์สำหรับวิเคราะห์สภาวะสุขภาพของมนุษย์ สามารถตรวจสอบความถูกต้องของผลลัพธ์กับงานวิจัยที่เกี่ยวข้องได้ [5][59][60] ปรากฏว่าเวิร์คโฟลว์สามารถค้นหาข้อมูลที่เกี่ยวข้องกับสภาวะต่างๆที่สนใจได้ถูกต้อง และความถูกต้องของผลลัพธ์การทำงานได้รับการรับรองจากนักชีวสารสนเทศจากศูนย์วิจัยเภสัชพันธุศาสตร์และสารสนเทศ ของโรงพยาบาลรามาธิบดี มหาวิทยาลัยมหิดล และศูนย์ความเป็นเลิศด้านชีววิทยาศาสตร์ของประเทศไทย [62]

6.3 ข้อเสนอแนะ

จากผลการทดลองเวิร์คโฟลว์สำหรับวิเคราะห์สภาวะสุขภาพของมนุษย์ บ่งบอกว่าจำเป็นต้องมีการปรับแต่ง หรือออกแบบเวิร์คโฟลว์ใหม่บ้างเล็กน้อย นั่นคือการเอาบริการที่ล้าสมัยออกและเพิ่มบริการใหม่เข้ามาแทนที่และจัดการเชื่อมโยงพอร์ตของข้อมูลต่างๆตามความต้องการ ในกรณีที่พบว่าแม้บริการจะมีหลายๆฟิลเตอร์ หากแต่มีฟิลเตอร์เพียงตัวเดียวที่ล้าสมัยไปแล้วก็จะทำให้ Data flow มีความผิดพลาดอยู่ เนื่องจากบริการจะต้องนำฟิลเตอร์ทั้งหมดเป็นไปเงื่อนไขของการคิวรี ดังนั้นฟิลเตอร์ทุกตัวจะต้องผ่านเงื่อนไขนั้น (AND condition) ซึ่งปัญหานี้จำเป็นต้องเปลี่ยนเป็นบริการใหม่เท่านั้นจึงจะสามารถทำงานต่อได้ ในกรณีของบริการที่มีเพียงแอตทริบิวต์ก็เช่นเดียวกัน หากเพียงแอตทริบิวต์เดียวที่ล้าสมัยไปแล้วก็ส่งผลให้ Data flow มีความผิดพลาดถึงแม้ว่าแอตทริบิวต์ตัวนั้นจะไม่ได้นำไปใช้ประมวลผลกับกระบวนการถัดไปก็ตาม

ฟิลเตอร์และแอตทริบิวต์ที่ล้าสมัย อันเนื่องมาจากการปรับปรุงของบริการไบโอมาร์ทเหล่านี้ จะไม่แสดงตัวที่หน้าต่างการปรับแต่งของบริการในโปรแกรมทาวเวอร์นา ดังนั้นผู้ใช้งานจึงไม่ทราบความล้าสมัย และไม่สามารถยกเลิกการเลือกฟิลเตอร์และแอตทริบิวต์เหล่านี้ได้ ผู้ใช้ที่ชำนาญอาจจะแก้เอกสาร Scum ของเวิร์คโฟลว์ด้วยตัวเองได้ แต่วิธีนี้ก็สร้างความสับสนและทำให้เกิดข้อผิดพลาดได้ง่าย

ในที่สุดแล้ว คำแนะนำก็คือการเพิ่มบริการใหม่ โดยยังเป็นชื่อบริการเดิมอยู่ซึ่งปรากฏอยู่ในรายงานการตรวจสอบเวิร์คโฟลว์ของระบบตรวจสอบ ที่นำเสนอในวิทยานิพนธ์นี้ ผู้ใช้งานสามารถเลือกบริการนี้ใหม่อีกครั้งจากหน้าต่าง *Service available panel* ในโปรแกรมทาวเวอร์นา จากนั้นเชื่อมต่อกับบริการและเชื่อมโยงข้อมูลที่เกี่ยวข้องเข้าด้วยกัน และในกรณีที่

บริการใหม่ que เพิ่มเข้ามานี้ไม่มีฟิลเตอร์หรือแอ็คตริบิวต์ที่ต้องการ จึงมีความจำเป็นต้องเลือกบริการอื่นๆที่อยู่ในหมวดหมู่เดียวกันมาใช้แทน

งานวิจัยในวิทยานิพนธ์นี้ สร้างเวิร์คโฟลว์สำหรับการตรวจสอบฟิลด์ข้อมูลเงื่อนไข และฟิลด์ข้อมูลผลลัพธ์ของบริการไป โอมาร์ทในสิ่งแวดล้อมของโปรแกรมทาวเวอร์นา เวิร์คโฟลว์สำหรับตรวจสอบบริการไป โอมาร์ทนี้ ช่วยในการปรับปรุงกระบวนการของเวิร์คโฟลว์ โดยไม่ต้องรอให้เวิร์คโฟลว์ทำงานเสร็จเสียก่อนที่จะทราบว่าเวิร์คโฟลว์ไม่สมบูรณ์ และไม่สามารถให้ผลลัพธ์ที่ต้องการได้ โดยสามารถประหยัดเวลาในการทำงานได้มากโดยไม่ต้องเสียไปกับการทำงานที่ไม่ได้ผลลัพธ์ใดๆ หรือได้ผลลัพธ์ที่ไม่ต้องการและสามารถลดเวลาในการดีบั๊กเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์ลงได้มาก นอกจากนี้เวิร์คโฟลว์สำหรับการตรวจสอบบริการไป โอมาร์ทยังนำเสนอข้อมูลที่เกี่ยวข้องกับบริการที่ทันสมัยนั้นๆสำหรับการตรวจสอบในอนาคตอีกด้วย

6.4 สรุป

ในบทนี้ ได้กล่าวถึงการทดสอบเวิร์คโฟลว์สำหรับการตรวจสอบความถูกต้องของบริการไป โอมาร์ท โดยใช้เวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์เป็นกรณีศึกษา รายงานผลการทดสอบเวิร์คโฟลว์ สามารถระบุฟิลเตอร์และแอ็คตริบิวต์ที่ล้าสมัยได้อย่างถูกต้อง หรือกล่าวอีกนัยหนึ่งก็คือ เวิร์คโฟลว์สำหรับการตรวจสอบสามารถระบุบริการไป โอมาร์ทที่ล้าสมัยได้อย่างถูกต้อง ซึ่งใช้เวลาน้อยกว่าครั้งหน้าที่และสามารถลดเวลาในการดีบั๊กเวิร์คโฟลว์ลงได้อย่างมีนัยสำคัญ รวมทั้งให้ข้อมูลที่สำคัญในการช่วยปรับแก้เวิร์คโฟลว์ได้ โดยการลบบริการเดิมและเพิ่มบริการใหม่ลงไป จากนั้นเชื่อมต่อการเชื่อมโยงข้อมูลกับบริการอื่นๆและทดสอบซ้ำ ซึ่งเวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์สามารถทำงานได้ถูกต้องโดยปราศจากข้อผิดพลาดใดๆ

บทที่ 7

สรุป

ในบทนี้จะกล่าวสรุปผลการวิจัยที่ได้ดำเนินการสำหรับวิทยานิพนธ์นี้ รวมทั้งข้อเสนอแนะต่างๆ เพื่อเป็นประโยชน์ต่อการทำวิจัยเกี่ยวกับเทคโนโลยีมายกริดและโปรแกรมทาวเวอร์นาต่อไป

7.1 ภาพรวมของงานวิทยานิพนธ์

วิทยานิพนธ์นี้เป็นการประยุกต์ใช้เทคโนโลยีมายกริด โปรแกรมทาวเวอร์นา และเว็บเซอร์วิสเพื่อแก้ไขปัญหาในงานวิจัยทางด้านชีวสารสนเทศศาสตร์ ซึ่งเป็นงานวิจัยที่ทำร่วมกับหน่วยงานวิจัยที่มีความชำนาญทางด้านชีวสารสนเทศศาสตร์ของมหาวิทยาลัยสงขลานครินทร์ และโรงพยาบาลรามาราชดี คณะแพทยศาสตร์ มหาวิทยาลัยมหิดล วิทยานิพนธ์นี้ได้บรรลุวัตถุประสงค์ดังต่อไปนี้

1) ได้นำเสนอวิธีการพัฒนาเว็รค์โฟลว์ และสร้างการเชื่อมต่อในเว็รค์โฟลว์สำหรับแก้ไขปัญหาในงานวิจัย การทำนายผลกระทบของสปีชของมนุษย์และการวิเคราะห์สปีชของกิ้ง โดยใช้เทคโนโลยีมายกริด, โปรแกรมทาวเวอร์นา, เว็บเซอร์วิสและใช้เว็รค์โฟลว์ทั้งสองนี้เป็นกรณีศึกษาของงานวิทยานิพนธ์

2) ได้นำเสนอกลไกในการลดเวลาในการพัฒนา และทดสอบเว็รค์โฟลว์โดยพัฒนาการตรวจสอบความถูกต้องของเว็รค์โฟลว์ก่อนการทำงานจริง โดยจะมุ่งเน้นไปที่การตรวจสอบบริการไบโอมาร์ทภายใต้สิ่งแวดล้อมของโปรแกรมทาวเวอร์นา ที่สามารถตรวจสอบคุณสมบัติของบริการไบโอมาร์ทได้ว่าทันสมัยพร้อมทำงานหรือล้าสมัยไปแล้ว ซึ่งสามารถระบุบริการที่ล้าสมัยไปแล้วในเว็รค์โฟลว์ และให้ข้อมูลที่มีนัยสำคัญต่อผู้ใช้งานจะได้มีความสะดวกและง่ายต่อการปรับแก้เว็รค์โฟลว์ซึ่งสามารถลดเวลาที่เสียไปได้อย่างมาก

3) ได้นำเสนอกฎในการลดปัญหาการเกินเวลา (Time-out) และความไม่คงเส้นคงวาในการทำงาน (Inconsistency) ที่ส่งผลให้เวิร์คโฟลว์ทำงานไม่สำเร็จโดยพัฒนาบริการท้องถิ่นขึ้นมาใช้งานเอง (Local web services)

วิทยานิพนธ์นี้บรรจุวัตถุประสงค์ทั้ง 3 ข้อดังกล่าว และจะสรุปผลลัพธ์ของแต่ละงานในวิทยานิพนธ์ดังนี้

7.2 การเลือกบริการที่เหมาะสมกับลักษณะของงานและการพัฒนาเวิร์คโฟลว์

งานวิทยานิพนธ์นี้ได้นำเสนอการเลือกบริการที่เหมาะสม และการวาดเวิร์คโฟลว์ โดยให้ข้อเสนอแนะที่เป็นกระบวนการทำงานแก่ผู้ใช้งาน โดยแบ่งวิธีการค้นหาและการเลือกบริการออกเป็น 2 วิธีคือ

- การค้นหาบริการจากเว็บพอร์ทัลของมายกริดซึ่งเรียกว่า Biological Web Service โดยสามารถเข้าถึงได้ที่ URL ดังนี้ <http://www.mygrid.org.uk/wiki/Mygrid/BiologicalWebServices>
- การค้นหาบริการจากอินเทอร์เน็ตโดยปลั๊กอินของโปรแกรมทาเวอร์น่าเอง หรือ feta ซึ่งเป็นการค้นหาเว็บเซอร์วิสโดยใช้พื้นฐานจากคำอธิบายของหน้าที่ในการทำงานของบริการนั้นๆ (Description) ด้วย Service ontology ซึ่งจะเก็บอยู่ที่รีจิสทรีของบริการหรือ UDDI ของมายกริด หากค้นหาบริการที่ต้องการได้ก็สามารถทดสอบการใช้งานจากโปรแกรมทาเวอร์น่าได้ทันที

นอกจากนี้ยังนำเสนอวิธีการพัฒนาหรือการวาดเวิร์คโฟลว์ โดยเริ่มต้นที่สมมุติฐานที่ว่าผู้วิจัยได้การออกแบบ Flowchart หรือ Data flow diagram ของระบบงานวิจัยเสร็จสิ้นแล้ว ผู้วิจัยมองเห็นอุปสรรค การไหลของข้อมูลและเอาท์พุทโดยรวมจากการออกแบบได้ โดยพร้อมที่จะพัฒนาเวิร์คโฟลว์ของระบบงานตามที่ได้ออกแบบ

7.3 บริการท้องถิ่นในการทำนายโครงสร้างต้นสายวิวัฒนาการของกุ้ง

ในการสร้างเวิร์คโฟลว์เพื่อวิเคราะห์สปีชีส์ในกุ้ง ซึ่งเป็นกรณีศึกษาเบื้องต้น วิทยานิพนธ์นี้ แต่เวิร์คโฟลว์ที่พัฒนาครั้งแรกไม่สามารถทำงานได้สำเร็จเนื่องจากเกิดปัญหาการเกิน

เวลาและความไม่คงเส้นคงวาในการทำงานระหว่างการถ่ายโอนข้อมูลระหว่างเว็บเซอร์วิสของ EBI ในส่วนของการสร้างโครงสร้างต้นไม้สายวิวัฒนาการของกิ้ง

วิทยานิพนธ์นี้ได้พัฒนาเว็บเซอร์วิสท้องถิ่นหรือบริการท้องถิ่น (Local web services) ขึ้นมาใช้งานเองภายในระบบเวิร์คโฟลว์ โดยใช้ Soaplab Analysis Tool เป็นเครื่องมือช่วยสร้างตัวเชื่อมประสานหรือ Web service interfaces ให้กับโปรแกรมประยุกต์ต่างๆที่เกิดปัญหาการกินเวลา เพื่อให้โปรแกรมทาวเวอร์น่าเข้าถึงได้ และสร้าง Java Web Service (JWS) ในการจัดการอินพุตเอาต์พุตภายในเวิร์คโฟลว์ด้วย บริการท้องถิ่นที่สร้างขึ้นสามารถแก้ไขปัญหาการกินเวลาและความไม่คงเส้นคงวาในการทำงานได้เป็นอย่างดี สามารถลดเวลาการทำงานลงได้อย่างมีนัยสำคัญเมื่อเปรียบเทียบกับการทำงานด้วยวิธีการแบบเดิม (คัดลอกแล้ววาง) และผลลัพธ์ที่ได้จากการทำงานของเวิร์คโฟลว์เพื่อวิเคราะห์สนิปในกิ้งนี้ ได้รับการตรวจสอบความถูกต้องจากศูนย์วิจัยจีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์

7.4 เวิร์คโฟลว์สำหรับวิเคราะห์สนิปของมนุษย์

วิทยานิพนธ์นี้ได้กล่าวถึงการพัฒนาวิธีการทำนายผลกระทบ อันเนื่องมาจากการเปลี่ยนแปลงสนิปต่างๆของมนุษย์ รวมทั้งสร้างระบบที่มีความสามารถในการคัดเลือกตำแหน่งสนิปที่น่าสนใจในการนำไปศึกษาต่อ ทางด้านหน้าที่การทำงานของสนิปต่างๆต่อไปด้วยเวิร์คโฟลว์ การแก้ปัญหาทางงานวิจัยใช้โปรแกรมทาวเวอร์น่า ที่มีความสามารถในการสร้างและจัดการรูปแบบข้อมูล รวมทั้งมีความสามารถในการดึงข้อมูลจากแหล่งต่างๆ โดยประสานงานกับเว็บเซอร์วิสชนิดต่างๆและแสดงผลออกมาในรูปแบบที่ง่ายต่อการนำไปศึกษา หรือใช้งานต่อได้และได้รับการตรวจสอบจากนักชีวสารสนเทศแล้วว่าทำงานได้ถูกต้องตรงตามความต้องการ และเร็วกว่าการทำงานตามปกติ และใช้เวิร์คโฟลว์นี้เป็นกรณีศึกษาหลักของวิทยานิพนธ์

7.5 เวิร์คโฟลว์สำหรับการตรวจสอบบริการไบโอมาร์ท

วิทยานิพนธ์นี้นำเสนอแนวคิดในการออกแบบ การพัฒนา และการทดสอบเวิร์คโฟลว์สำหรับตรวจสอบการเปลี่ยนแปลง ของบริการไบโอมาร์ทในเวิร์คโฟลว์วิเคราะห์ของมนุษย์ของงานวิจัยเภสัชพันธุศาสตร์แบบอัตโนมัติ และเวิร์คโฟลว์ที่พัฒนาขึ้นใช้บริการท้องถิ่นที่

โปรแกรมทาวเวอร์นาเตรียมไว้ให้สำหรับการสร้างสคริปต์ของภาษา Java ขึ้นมาทำงานตาม อัลกอริทึมที่ออกแบบไว้ จากการทดสอบด้วยเวิร์คโพล์กรณีศึกษาอย่างง่ายพบว่า เวิร์คโพล์ สำหรับตรวจสอบการเปลี่ยนแปลงของบริการไปโอมาร์ทซึ่งสามารถระบุได้ว่า บริการไปโอมาร์ท นั้นๆ ได้ล้าสมัยไปแล้วได้อย่างถูกต้อง

จากนั้นได้นำเสนอผลการทดสอบของเวิร์คโพล์ สำหรับการตรวจสอบความ ถูกต้องของบริการไปโอมาร์ท โดยมีเวิร์คโพล์วิเคราะห์สัณนิษของมนุษย์เป็นกรณีศึกษา ซึ่งสามารถ ระบุบริการไปโอมาร์ทที่ล้าสมัยได้อย่างถูกต้องโดยใช้เวลาน้อยกว่าครึ่งนาที และรายงานผลการ ตรวจสอบยังสามารถลดเวลาในการดีบั๊กเวิร์คโพล์ลงได้อย่างมีนัยสำคัญ ซึ่งให้ข้อมูลที่สำคัญใน การช่วยปรับแก้เวิร์คโพล์ได้ และทดสอบซ้ำหลังจากปรับแก้แล้ว เวิร์คโพล์กรณีศึกษาสามารถ ทำงานได้ถูกต้องโดยปราศจากข้อผิดพลาดใดๆ และความถูกต้องของผลลัพธ์การทำงานได้รับการ รับรองจากนักชีวสารสนเทศ จากศูนย์วิจัยเภสัชพันธุศาสตร์และสารสนเทศของโรงพยาบาล ราชามาธิบดี คณะแพทยศาสตร์ มหาวิทยาลัยมหิดล และศูนย์ความเป็นเลิศด้านชีววิทยาศาสตร์ของ ประเทศไทย

7.6 อุปสรรคและปัญหา

- การพัฒนาระบบงานแบบเวิร์คโพล์ในโปรแกรมทาวเวอร์นา ยังถือว่าเป็น งานวิจัยเฉพาะด้านและต้องเป็นไปตามมาตรฐานที่โครงการมายกริดนำเสนอ เช่น สถาปัตยกรรม ของ Soaplab หรือภาษา Scufi ของเอกสารเวิร์คโพล์ ซึ่งผู้วิจัยต้องทำความรู้จักกับมาตรฐานใหม่นี้ ทั้งนี้โครงการมายกริดก็ไม่ได้สนับสนุนการใช้ภาษา Scufi ภายนอกสิ่งแวดล้อมของโปรแกรม ทาวเวอร์นา ภาษา Scufi เป็นภาษาที่เฉพาะเจาะจงสำหรับระบบงานแบบเวิร์คโพล์และประยุกต์ใช้ กับมาตรฐานอื่นๆ ได้ไม่ถนัดนัก

- โปรแกรมทาวเวอร์นาอนุญาตให้ผู้วิจัยสามารถสร้างอัลกอริทึม สำหรับการ ทำงานได้ตามความต้องการ บริการท้องถิ่นบนเซลล์ที่โปรแกรมนำเสนอต่อผู้วิจัยนั้นรองรับการ โปรแกรมภาษา Java ดังนั้นผู้วิจัยจะต้องเขียนภาษาสคริปต์ของภาษา Java เท่านั้น ปัญหาที่เกิดขึ้น คือ หากอัลกอริทึมมีกระบวนการทำงานที่ซับซ้อน จะทำให้การดีบั๊กความถูกต้องของสคริปต์ทำได้ ยาก เนื่องจากจะต้องสั่งให้เวิร์คโพล์ทำงานจริงๆ จึงจะสามารถดีบั๊กความถูกต้องของสคริปต์ได้

7.7 ข้อเสนอแนะ

- การพัฒนาเว็บเซอร์วิสท้องถิ่นหรือบริการท้องถิ่น (Local web services) ขึ้นมาใช้งานเองภายในระบบเวิร์คโพล์เพื่อวิเคราะห์สปีโนกิ้ง เป็นบริการที่ทำงานบนโครงสร้างพื้นฐานของระบบคอมพิวเตอร์คลัสเตอร์ แต่เนื่องจากชุดแอ็พพลิเคชัน Emboss และ PHYLIB ที่ทางศูนย์จีโนมิกส์และชีวสารสนเทศแห่งมหาวิทยาลัยสงขลานครินทร์มีใช้งาน เป็น โปรแกรมในการทำนายโครงสร้างสายวิวัฒนาการยังคงเป็นแอ็พพลิเคชันแบบ Sequential ดังนั้นบริการท้องถิ่นที่ได้พัฒนาในวิทยานิพนธ์นี้ จึงยังไม่สามารถดึงความสามารถด้านการประมวลผลแบบขนานของระบบคอมพิวเตอร์คลัสเตอร์มาใช้งานได้เต็มที่ บริการท้องถิ่นที่พัฒนาทำงานเพียงโหนดเดียวเท่านั้นคือที่โหนด front end ของระบบคลัสเตอร์ Hanuman ของศูนย์กริคมมหาวิทยาลัยสงขลานครินทร์ ข้อเสนอแนะคือ นอกจากจะมองว่าบริการแต่ละบริการในเวิร์คโพล์จะทำงานเป็นลักษณะการกระจายแล้ว บริการท้องถิ่นที่พัฒนาก็อาจจะปรับปรุงให้ทำงานแบบขนานได้ ซึ่งต้องใช้ชุดแอ็พพลิเคชัน Emboss และ PHYLIB ที่ทำงานแบบขนานสำหรับงานบนระบบคอมพิวเตอร์คลัสเตอร์ ซึ่ง ณ ปัจจุบันที่ทำวิทยานิพนธ์นี้ยังไม่มีชุดแอ็พพลิเคชัน Emboss และ PHYLIB ที่ทำงานแบบขนานสำหรับงานบนระบบคอมพิวเตอร์คลัสเตอร์ออกมาให้ใช้งาน อย่างไรก็ตาม ยังมีงานวิจัยที่พัฒนาระบบการกระจายงานในระบบขนานสำหรับชุดแอ็พพลิเคชัน Emboss ขึ้นมาใช้งานเอง [63] นอกจากนี้สถาบัน National Institutes of Health (NIH) ซึ่งเป็นผู้ให้บริการระบบคอมพิวเตอร์คลัสเตอร์ Biowulf [64] พัฒนาโปรแกรมชื่อ Swarm [65] สำหรับช่วยกระจายงานให้เป็นแบบขนานบนระบบคอมพิวเตอร์คลัสเตอร์ได้ และนักวิจัยที่สนใจสามารถดาวน์โหลดโปรแกรม Swarm มาใช้งานกับระบบคอมพิวเตอร์คลัสเตอร์ของตนเองได้อีกด้วย แต่ทั้งนี้การทำงานบนระบบคอมพิวเตอร์คลัสเตอร์ ซึ่งมีโครงสร้างพื้นฐานที่เป็นฮาร์ดแวร์ที่ทรงประสิทธิภาพย่อมทำให้บริการท้องถิ่นมีความน่าเชื่อถือในการทำงานและมีความสามารถในเข้าถึงสูง (Accessibility)

- การเสนอแนวคิดในการออกแบบ การพัฒนา และการทดสอบเวิร์คโฟลว์สำหรับตรวจสอบการเปลี่ยนแปลงของบริการ ที่ใช้ในสำหรับวิเคราะห์สลิปของมนุษย์แบบอัตโนมัติ เป็นการตรวจสอบกับบริการไปโอมาร์ทเท่านั้น เนื่องจากเป็นบริการหลักที่ใช้ในเวิร์คโฟลว์กรณีศึกษา ข้อเสนอแนะคือ การพัฒนาระบบตรวจสอบนี้ให้สามารถตรวจสอบกับบริการอื่นๆได้ด้วย เช่น BioMoby หรือ e-Utilities ของ NCBI เป็นต้น ในขั้นตอนการปรับแก้เวิร์คโฟลว์หลังการตรวจสอบก็อาจจะนำเสนอกลไกในการแก้ไขเวิร์คโฟลว์แบบอัตโนมัติได้ นอกนี้กลไกตรวจสอบเหล่านี้ ก็อาจจะไม่จำเป็นต้องอยู่ในรูปแบบของเวิร์คโฟลว์เสมอไปก็ได้ เช่น อาจจะอยู่ในรูปแบบของเว็บเซอร์วิส เป็นต้น

เอกสารอ้างอิง

- [1] Oinn, T., Addis, M., Ferris, J., Marvin, D., Senger, M., Greenwood, M., Carver, T., Glover, K., Pocock, M.R., Wipat, A. and Li, P., “Taverna: a tool for the composition and enactment of bioinformatics workflows,” *BIOINFORMATICS*, 2004 Vol. 20, no. 17, pp. 3045–3054, doi:10.1093/bioinformatics/bth361.
- [2] Wolstencroft, K.; Oinn, T.; Goble, C.; Ferris, J.; Wroe, C.; Lord, P.; Glover, K.; Stevens, R., “Panoply of utilities in Taverna,” *First International Conference on e-Science and Grid Computing*, 5-8 Dec. 2005.
- [3] Hull, D., Wolstencroft, K., Stevens, R., Goble, C., Pocock, M.R., Peter Li, P. and Oinn, T., “Taverna: a tool for build and running workflows of services,” *Nucleic Acids Research*, 2006, Vol. 34, Web Server issue W729–W732.
- [4] Turi, D., Missier, P., Goble, C., De Roure, D., Oinn, T., “Taverna Workflows: Syntax and Semantics, *e-Science and Grid Computing*,” *IEEE International Conference on* 10-13 Dec. 2007, pp. 441–448.
- [5] Wellcome Trust Case Control Consortium, “Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls,” *Nature* 2007, 447, pp. 661-678.
- [6] SNP. 2009. <http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp>.
- [7] Young C. Song, Edward Kawas, Ben M. Good , Mark D. Wilkinson and Scott J. Tebbutt, “DataBiNS: a BioMoby-based data-mining workflow for biological pathways and non-synonymous SNPs,” *Bioinformatics Advance Access* originally published online on January 18, 2007, *Bioinformatics* 2007 23(6), pp. 780-782, doi:10.1093/bioinformatics/btl648.
- [8] Taverna project website. 2009. <http://taverna.sourceforge.net>.
- [9] myGrid. 2009. <http://www.mygrid.org.uk/>.
- [10] R. Stevens, H.J. Tipney, C. Wroe, T. Oinn, M. Senger, P. Lord, C.A. Goble, A. Brass and M. Tassabehji, “Exploring Williams-Beuren Syndrome Using myGrid in *Proceedings of 12th International Conference on Intelligent Systems in Molecular Biology*,” 31st Jul-4th Aug 2004, Glasgow, UK, published *Bioinformatics* Vol. 20 Suppl. 1 2004, i303-i310.
- [11] BioMart. 2009. <http://www.biomart.org/>.

- [12]Damian Smedley, Syed Haider, Benoit Ballester, Richard Holland, Darin London, Gudmundur Thorisson and Arek Kasprzyk., "BioMart - biological queries made easy," BMC Genomics 2009, 10:22doi:10.1186/1471-2164-10-22, Published 14 January 2009, <http://www.biomedcentral.com/1471-2164/10/22>.
- [13]Wanna, W., Rolland, J.L., Bonhomme, F. and Phongdara, A., "Population genetic structure of *Penaeus merguensis* in Thailand based on nuclear DNA variation," J. Exp. Mar. Bio, 2004, Ecol. 311, pp. 63-78.
- [14]ธีราวุฒิ ภู่อันติสัมพันธ์. 2549. การศึกษา Single nucleotide polymorphisms ของยีน amylase ใน กุ้งแชบ๊วย. วิทยานิพนธ์วิทยาศาสตรบัณฑิต, สาขาวิชาชีวเคมี คณะวิทยาศาสตร์, มหาวิทยาลัยสงขลานครินทร์.
- [15]Phylogenetic trees. 2009. [http:// cnx .org /content /m11052 /latest/](http://cnx.org/content/m11052/latest/).
- [16]BioDAS. 2009. http://www.biodas.org/wiki/Main_Page.
- [17]Dasty2, a protein DAS client. 2009. <http://www.ebi.ac.uk/dasty/>.
- [18]MySQL. 2009. www.mysql.com/.
- [19]Oracle. 2009. <http://www.oracle.com/index.html>.
- [20]PostgreSQL. 2009. www.postgresql.org/.
- [21]GNU Lesser General Public License, GNU Project - Free Software. 2009. <http://www.gnu.org/licenses/lgpl.html>.
- [22]Peter M. Rice, Alan J. Bleasby, Syed A. Haider, Jon C. Ison, Shaun McGlinchey, Mahmut Uludag, "EMBRACE: Bioinformatics Data and Analysis Tool Services for e-Science," e-science, pp.146, Second IEEE International Conference on e-Science and Grid Computing (e-Science'06), 2006, DOI Bookmark: <http://doi.ieeecomputersociety.org/10.1109/E-SCIENCE.2006.57>.
- [23]Java Database Connectivity (JDBC), Java SE Technology, Database. 2009. <http://java.sun.com/javase/technologies/database/>
- [24]Nevirapine, FDA Public Health Advisory for Nevirapine (Viramune). 2009. <http://www.fda.gov/CDER/drug/advisory/nevirapine.htm>.
- [25]Gold, S., "Stevens-Johnson Syndrome," Br Med J. 1953 May 16; 1(4819): 1111. PMID: PMC2016447. Available. 2009. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2016447>.

- [26]Lude Franke, Harm van Bakel, Like Fokkens, Edwin D. de Jong, Michael Egmont-Petersen, Cisca Wijmenga, "Reconstruction of a functional human gene network, with an application for prioritizing positional candidate genes." *Am J Hum Genet* 2006 Jun; 78(6):1011-25.
- [27]dbSNP Overview, NCBI. 2009. http://www.ncbi.nlm.nih.gov/projects/SNP/get_html.cgi?whichHtml=overview.
- [28]KEGG: Kyoto Encyclopedia of Genes and Genomes. 2009. <http://www.genome.jp/kegg/>.
- [29]Lei Bao, Mi Zhou and Yan Cui, "nsSNPAnalyzer: identifying disease-associated nonsynonymous single nucleotide polymorphisms," *Nucleic Acids Research* 2005 33(Web Server Issue):W480-W482; doi:10.1093/nar/gki372.
- [30]Karchin, R., Diekhans, M., Kelly, L., Thomas, DJ., Pieper, U, Eswar, N., Haussler, D. and Sali, A., "LS-SNP: large-scale annotation of coding non-synonymous SNPs based on multiple information sources," *Bioinformatics*. 2005 Jun 15;21(12):2814-20. Epub 2005 Apr 12.
- [31]ศิวคณ ไชยศิริ, ศิริชัย โจรณ์วิภาค และผศ.ดร. ภูชงค์ อุทโยภาส, กริดเซอร์วิส – เว็บเซอร์วิสสายพันธุ์ใหม่. 2009. <http://apstf.cpe.ku.ac.th/oldweb/download/GridService.pdf>.
- [32]Web Service articles. 2009. <http://www.service-architecture.com/web-services/articles/>.
- [33]SOAP Specifications. 2009. <http://www.w3.org/TR/soap/>.
- [34]Web Service Definition Language (WSDL). 2009. www.w3.org/TR/wsdl.
- [35]PSU-Grid. 2009. <http://psu-grid.coe.psu.ac.th/>.
- [36]FASTA format description. 2009. <http://www.ncbi.nlm.nih.gov/blast/fasta.shtml>.
- [37]EBI Web Service for EMBOSS-4.1 application. 2009. <http://www.ebi.ac.uk/soaplab/emboss4/services>.
- [38]Damkliang, K., Tandayya, P., Sangket, U., and Phusantisampan, T. 2009. "Workflow and Supporting Services for Single Nucleotide Polymorphisms Analysis Gene Using Taverna." *Proceedings of 2009 International Conference on Information Management and Engineering (ICIME 2009)*, Kuala Lumpur, Malaysia, April 3-5, 2009. pp. 27-31.
- [39]European Bioinformatics Institute. 2009. www.ebi.ac.uk/
- [40]Web API for Bioinformatics. 2009. <http://xml.ddbj.nig.ac.jp/index.html>
- [41]NCBI BLAST. 2009. www.ncbi.nlm.nih.gov/BLAST/
- [42]EBI Tools, InterProScan. 2009. www.ebi.ac.uk/Tools/InterProScan/

- [43]SeqVISTA: Towards an Integrative Platform for BioSeqs Analyses. 2009. <http://zlab.bu.edu/SeqVISTA/>.
- [44]TreeView X. 2009. <http://darwin.zoology.gla.ac.uk/~rpage/treeviewx/index.html>.
- [45]Web Services - Axis. 2009. <http://ws.apache.org/axis/>.
- [46]The Gene Ontology. 2009. www.geneontology.org/.
- [47]GNU Lesser General Public License, GNU Project, Free Software Foundation (FSF). 2009. <http://www.gnu.org/licenses/lgpl.html>.
- [48]PHYLIP Home Page. 2009. <http://evolution.genetics.washington.edu/phylip.html>.
- [49]AJAX Command Definition. 2009. <http://embooss.sourceforge.net/developers/acd/>.
- [50]Ensembl Genome Browser. 2009. <http://www.ensembl.org/index.html>.
- [51]Homo sapiens-UniGene. 2009. <http://www.ncbi.nlm.nih.gov/UniGene/UGOrg.cgi?TAXID=9606>.
- [52]Haman (Homo sapiens). 2009. http://ensembl.genomics.org.cn:8050/Homo_sapiens/index.html.
- [53]Entrez cross-database search. 2009. <http://www.ncbi.nlm.nih.gov/sites/gquery>.
- [54]OMIM Home. 2009. <http://www.ncbi.nlm.nih.gov/sites/entrez?db=omim>.
- [55]MOTIF: Searching Protein and Nucleic Acid Sequences Motifs. 2009. <http://motif.genome.jp/>.
- [56]KEGG: Kyoto Encyclopedia of Gene and Genomes. 2009. <http://www.genome.jp/kegg/>.
- [57]KEGG Pathway Database. 2009. <http://www.genome.ad.jp/kegg/pathway.html>.
- [58]XML Path Language (XPath). 2009. <http://www.w3.org/TR/xpath>.
- [59]Han J, Kraft P, Nan H, Guo Q, Chen C, et al. (2008)., "A Genome-Wide Association Study Identifies Novel Alleles Associated with Hair Color and Skin Pigmentation," *PLoS Genet* 4(5): e1000074. doi:10.1371/journal.pgen.1000074.
- [60]Damkhang, K., Tandayya, P., Pasomsuba, E., Chantratita, W., and Mahasirimongkol, S. 2009. "Taverna Workflow for Validating Biomart Services." *Proceedings of 2009 International Conference on Computer Design and Applications (ICCD 2009)*, Singapore, May 15-17, 2009. pp. 459-463.
- [61]Axel M Hillmer, Felix F Brockschmidt, Sandra Hanneken, Sibylle Eigelshoven, Michael Steffens, Antonia Flaquer, Stefan Herms, Tim Becker, Anne-Katrin Kortüm, Dale R Nyholt,

Zhen Zhen Zhao, Grant W Montgomery, Nicholas G Martin, Thomas W Mühleisen, Margrieta A Alblas, Susanne Moebus, Karl-Heinz Jöckel, Martina Bröcker-Preuss, Raimund Erbel, Roman Reinartz, Regina C Betz, Sven Cichon, Peter Propping, Max P Baur, Thomas F Wienker, Roland Kruse and Markus M Nöthen, "Susceptibility variants for male-pattern baldness on chromosome 20p11," *Nature Genetics* 40, 1279-1281 (2008). Published online: 12 October 2008 | doi:10.1038/ng.228

[62] ศูนย์วิจัยเภสัชพันธุศาสตร์และสารสนเทศ. 2009. <http://www.pharmagtc.net/>.

[63] K. Podesta, M. Crane, H.J. Ruskin, "A sequence-focused parallelisation of EMBOSS on a cluster of workstations," *International Conference on Computational Science & its Applications (ICCSA '04)*, Assisi, Italy, May 2004, Publication in Springer Lecture Notes in Computer Science.

[64] Biowulf at the NIH. National Institutes of Health (NIH). 2009. <http://biowulf.nih.gov/>.

[65] Swarm on the Biowulf Linux Cluster. National Institutes of Health (NIH). 2009. <http://biowulf.nih.gov/apps/swarm.html>.

ภาคผนวก ก

บทความวิชาการระดับนานาชาติชิ้นที่ 1

เรื่อง

**Taverna Workflow and Supporting Services for
Single Nucleotide Polymorphisms Analysis**

งานประชุม

**2009 International Conference on Information Management and Engineering
(ICIME 2009), Kuala Lumpur, Malaysia, 3rd - 5th April 2009.**

PP. 27-31. ISBN: 978-1-4244-3774-0.

Taverna Workflow and Supporting Services for Single Nucleotide Polymorphisms Analysis

Kasikrit Damkliang¹ and Pichaya Tandayya²
Department of Computer Engineering,
Faculty of Engineering, Prince of Songkla University,
Hat Yai, Songkhla Thailand
s5010120123@psu.ac.th¹, pichaya@coe.psu.ac.th²

Theerawut Phusantisampan³ and Unitsa Sangket⁴
The Centre for Genomics and Bioinformatics Research,
Prince of Songkla University,
Hat Yai, Songkhla, Thailand
5110230005@psu.ac.th³, unitsa.s@psu.ac.th⁴

Abstract—This paper describes a Taverna workflow environment with supporting web services for Single Nucleotide Polymorphisms (SNP) analysis of genes in shrimp. The workflow is used to investigate SNPs advantages for gene development. Our first workflow environment did not give satisfying results for the phylogenetic tree processes because the data flows between the web services would frequently time-out. Our solution employs local web service interfaces in a cluster, which Taverna accesses with the aid of Soaplab analysis tools. The cluster enables the modified workflow to rapidly process the large amount of sequenced data required by phylogenetic trees, and is highly reliable.

Keywords—Taverna, Workflow, Web Service, SNP, Phylogenetics

I. INTRODUCTION

Bioinformatic researchers use to resort to copying and pasting between web pages to combine the results of multiple services or programs. This was labor intensive, error prone, and not scalable, due to standard differences, and research complexity [1, 3]. In response, workflow methodologies were developed to help compose the resources into a single logical system.

This paper presents a workflow environment based on Taverna [2, 4, 5] for Single Nucleotide Polymorphisms (SNP) [6] analysis of shrimp. We focus on the *Amylase* gene because it is a good model for studying the adaptation processes and phylogeny. *Penaeus merguensis* samples were collected from three of a multigene family regions of The Gulf of Thailand : Trad (TDE, N=3), Songkhla (SKE, N=3), and Surat Thani (SRE, N=3) [7]. Two or three independent clones were sequenced for each specimen.

A previous study [7] reported on the genetic variation of *P.merguensis* in Thailand using nuclear DNA markers (ACT1, D7MICRO and PvAmy). High intraspecific variation in PvAmy (FST= 0.324) was observed. This raises the question of the possibility of the PvAmy locus undergoing strong environmental selection, linked to variations in the ecological factors. Also, does the same phenomenon happen in other areas in Southeast Asia, such as the Malaysian peninsula. To examine this question, we study the polymorphism of the *Amy* gene using primer based on a conserved sequence from Exon7-10. This lets us

evaluate the level of genetic variation of *P.merguensis* from the Gulf of Thailand (Trad, Surat Thani and Songkhla). Moreover, SNPs may become very important in future studies of agriculture, such as growth and resistance to disease infection.

A major issue for a composed workflow is when it cannot finish its processing because data flows between the web services time-out in the phylogenetic tree [8] predicting process at the EBI [9]. This occurs when the processing involves for more than 2,000 bootstrapping replicates. It is possible to execute the predicting process locally but then we need to copy and paste data between the workflow, which is too much time-consuming.

A better solution is to build local web services that talk to the workflow, focusing on the building of the phylogenetic tree of nucleotide and Amino Acid sequences. The web service interfaces for each phylogenetic tree process in the workflow can then be accessed by Taverna. Our local web service interfaces were implemented with Soaplab [10] analysis tools.

Our workflows were implemented with BLAST (Basic Local Alignment and Search Tool) [11], multiple sequence alignment in nucleotide and Amino Acids, prediction of the protein secondary structure, and prediction of the nucleotide and Amino Acid phylogenetic tree [8]. We use Taverna to compose these web services into a workflow with services from XML Central of DDBJ [12], EBI [9], Java local services, local web services, and shim services. The correctness of each service was verified by The Centre for Genomics and Bioinformatics Research at Prince of Songkla University.

II. TAVERNA, WORKFLOW AND SOAPLAB INTRODUCTION

The myGrid e-Science Taverna tool is an open source workflow composer to orchestrate bioinformatic web services and existing applications into workflows, or to help compose local or distributed resources into a single logical system [5, 13].

Taverna can access over 30,000 services in the bioinformatics domain, and will grow to include other fields such as astronomy, chemoinformatics, health informatics in the future

Taverna can add service endpoints, such as WSDL documents, for improved portability, and can access other types of service, such as, local Java services, BioMoby and BioMart services, and Beanshell scripts [5, 13, 14].

Soaplab is a framework for exposing command-line tools as web services [10]. Soaplab Services containers manage web service jobs, such as Tomcat. It offers a various web services, and provides pointers to the AppLab Server which wraps existing command-line tools. Soaplab is a component of the myGrid project [13].

III. CASE STUDY BACKGROUND

A. Analysis of the Copying and Pasting method

SNP analysis used to rely on copying and pasting between web pages and off-line applications. The analysis processes include sequence analysis with NCBI BLAST [11], multiple alignments with ClustalX [15] and GeneDoc [16], translations of nucleotide sequences to Amino Acid sequences using Protein Sequence Analysis (PSA) [17], the study of evolution relationships by phylogenetic tree with BioEdit [18], and the prediction of secondary structure with PSIPRED [19]. The relationship between these processes is shown as a flowchart in Figure 1.

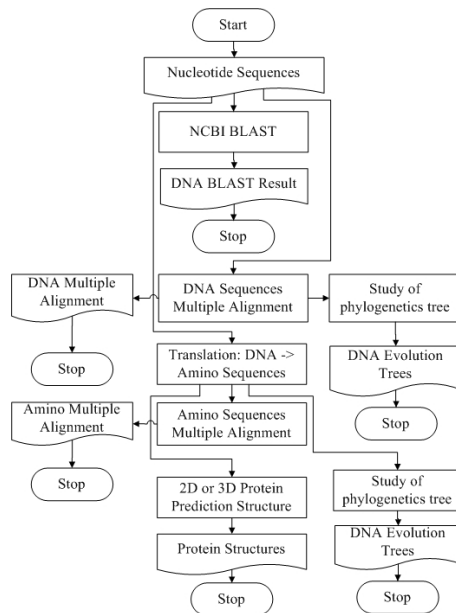


Figure 1. The Flowchart of the Shrimp SNP Analysis

The major web services used for making the workflow are Emboss application services at EBI [9]. Emboss offers many well-respected and broadly known web services.

The BLAST process is invoked by XML Central of DDBJ [20], using configurable databases such as DDBJNEW. BLAST is easy to access via Taverna or a web browser. The workflow also employs the Shim service, a small service that manages the data flow in and out of each service.

The workflow diagram is shown in Figure 2. Table I shows the Soaplab Web Services in the workflow, and the

string parameters configure each service [21] are listed in Table II.

In phylogenetics, an evolutionary tree produced with at least 1,000 bootstrapping replicates gives satisfying results because more replicates generate a plausible or correct tree [21]. We plan to utilize 2,000 replicates for precision results.

Figure 2 shows

- The use of existing services by Taverna.
- The prediction of a phylogenetic tree requires the 'seqboot' service or a bootstrapped sequences algorithm. 'seqboot' reads in a data set, and produces multiple data sets by sampling.
- The workflow input employs the shrimp nucleotide fasta format sequence [22] from a wet lab. There were to about 15-27 sequences; each consisting of 250 base pairs from the Exon7-10. An example of the fasta format sequences is shown below:

```
>SRE26
GTGGCGAAGCCATATCCAGCGCGGAGTATGTTGGCAACGGTCGTGTGACGG
AGTTCAGGTACGGCAAGTACCTGGGCGAGGCCTTCCGCGGCAACAACCCAGC
TGAATACTCTCAACAACCTTCGGCGAAGGTTGGGGCATGATTGACCGGCATG
ACGCCCTGGTCTTCATTGACAACCCAGCACAACCCAGAGAGGCCATGGTGCTG
AGGAGACATGATCCTTACTTCCGTGTCTCTAAGTGGTACAAGGA
>TDE8
GTGGCGAAGCCATATCCAGCGCGGAGTATGTTGGCAACGGTCGTGTGACGG
AGTTCAGGTACGGCAAGTACCTGGGCGAGGCCTTCCGCGGCAACAACCCAGC
TGAATACTCTCAACAACCTTCGGCGAAGGTTGGGGCATGATTGACCGGCATG
ACGCCCTGGTCTTCATTGACAACCCAGCACAACCCAGAGAGGCCATGGTGCTG
GAGGAGACATGACCCTTACTTCCGTGTCTCTAAGTGGTACAAGGA
>SKE42
GTGGCGAAGCCATATCCAGCGCGGAGTATGTTGGCAACGGTCGTGTGACGG
AGTTCAGGTACGGCAAGTACCTGGGCGAGGCCTTCCGCGGCAACAACCCAGC
TGAATACTCTCAACAACCTTCGGCGAAGGTTGGGGCATGATTGACCGGCATG
ACGCCCTGGTCTTCATTGACAACCCAGCACAACCCAGAGAGGCCATGGTGCTG
GAGGAGACATGATCCTTACTTCCGTGTCTCTAAGTGGTACAAGGA
```

- The process of predicting a phylogenetic utilizes 'seqboot', 'fdnapars' or 'fprotpars', and 'consense.'

The test results are shown in Table III. The workflow did not always produce results when using 2,000 bootstrapping replicates. The data flow between 'seqboot' → 'fdnapars' and 'seqboot' → 'fprotpars' sometimes terminates due to fetching time-out error or inconsistencies in retrieving the results.

B. Design the Solution

We solved the time-out connection problem by building local web services for predicting the phylogenetic tree. The services are hosted on Hanuman, a Rock cluster in the Grid Center at Prince of Songkla University [23].

The Soaplab Analysis Tool was used to implement the web service interfaces. Two types of web services were built, as shown in Figure 3. The application for predicting a phylogenetic tree was wrapped by Soaplab using the nucleotide acid and amino acid parsimony algorithm, I/O was managed using Java web service (JWS) [24].

C. Implementation

PHYLIB 3.6b [25] is an Emboss open source package for predicting a phylogenetic tree. PHYLIB is installed on Hanuman, employing the Rocks Linux Cluster Operating System version 4.3 [26] and Tomcat v.5.0.28 [27] as the web service container.

The key concept is to build an XML services description file for the Soaplab from an ACD [28] file list. An ACD file

must be written for each application program that acts as a web service [10].

TABLE I. EBI SOAPLAB WEB SERVICES USED IN THE WORKFLOW

Service Name	Responsibility
transeq	Translate nucleic acid Sequences to amino acid sequence
gamier	Predict protein secondary structure
fmma_[DNA/Prot]	Multiple alignment program interface to ClustalW program
prettyplot	Display aligned sequences, with coloring and boxing
fseqboot_[DNA/Prot]	Bootstrapped sequences algorithm
fdnapars	DNA parsimony algorithm
fprotpars	Protein parsimony algorithm
fconsense_[DNA/Prot]	Majority-rule and strict consensus tree

TABLE II. STRING PARAMETERS CONFIGURED IN THE WORKFLOW

Parameter Name	Responsibility	Value
Database	Database name used by blast for searching	DDBJNEW
Frame3	Translation using frame 3 [6]	3
osformat_[DNA/Prot]	Output sequences for phylogenetics	phylip
seqtype_[DNA/Prot]	Sequence type for bootstrapping	DNA = d Protein = p
Reps_[DNA/Prot]	Bootstrapping number for replication	Vary on test cases
rootTree_[DNA/Prot]	Produce the phylogenetics root tree	Yes

TABLE III. TESTING CASES AND RESULTS

No of Replicates	Time Used (min.)
100	02.50
250	12.20
500	20.29
1,000	20.40
2,000	Failed (time-out or inconsistency)

The script shown below is an example ACD file for 'fseqboot', 'fdnapars', 'fprotpars', and 'fconsense.' The ACD file configures the web service such as the location of an executable program or input file, the parameter '-auto' is for enabling the application to be non-interactive, the output file system location, or which web service is needed to retrieve the output from the host system

```
string: auto
[parameter: "Y" comment: "display false" default: "-auto" comment: "defaults"]
```

The Soaplab acd2xml tool was employed to translate the ACD files list into an XML file, Soaplab deployed the web services using the configuration parameters loaded from the file. The workflow for these local web services was composed in Taverna as shown in Figure 4. The new workflow with local web services is in the black frames.

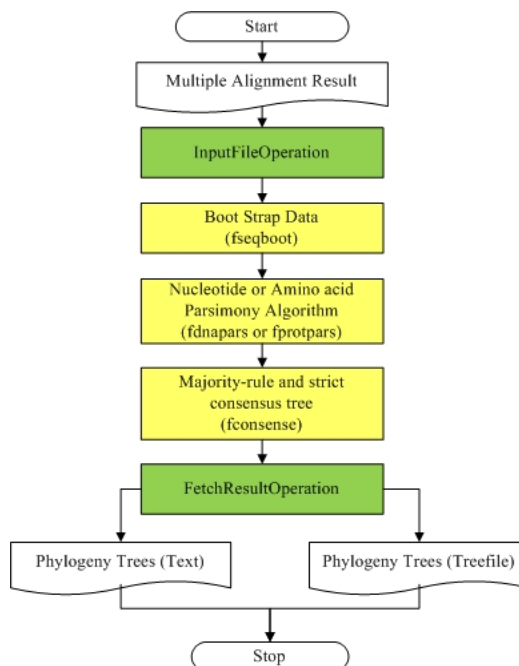


Figure 3. The Flowchart for Predicting the Phylogenetic Tree

IV. EXPERIMENTAL RESULTS

Table IV shows the results of our testing of the Figure 4 workflow, with the number of DNA bootstrapping replicates extended to 50,000. Table V compares manual processing (copying and pasting) versus automatic workflow processing. The first workflow failed due to a time-out error. With 2,000 DNA bootstrapping replicates, manual processing took 15 minutes whereas the automatic processing using the modified workflow took 2.24 minutes. The reduction in processing time is highly significant.

When we increased the number of DNA bootstrapping replicates to 20,000 and then to 50,000; the modified workflow took 15.08 minutes and 36.17 minutes respectively. We did not replicate these tests using manual processing because it is too difficult, or impossible, to manually copy and paste such large amount of data from and to web pages.

Fig. 5 and Fig. 6 show DNA and Amino Acid phylogenetic trees, verified by The Center for Genomics and Bioinformatics at Prince of Songkla University. The evolutionary relationships amongst the three different shrimp sources is quite clear, this knowledge will aid future studies of the selected species.

The DNA phylogenetic tree in Figure 5 highlights many genetic relationships. For example, TDE24 and SKE13 are in the same group of species, and have a closer relationship with SRE17 than with SRE16. The Amino Acid phylogenetic tree in Figure 6 shows that SKE50 and SKE42 are in same group of species, and were separated from TDE8.

TABLE IV. TESTING CASES AND RESULTS OF THE NEW WORKFLOW

No of Replicates	Time Used (Min.)
100	02.12
250	02.06
500	02.06
1,000	02.06
2,000	02.24
3,000	03.04
4,000	03.47
5,000	03.52
10,000	08.09
20,000	15.08
30,000	22.02
40,000	29.22
50,000	36.17

TABLE V. TESTING CASES AND RESULTS IN COMPARISON

No of Replicates	Test Case			Running Time (Min.)
	Manual	Automatic Workflow		
		No local service	With local web services	
2,000	X			15
2,000		X		Failed
2,000			X	2.24
20,000	X			Difficult or impossible to test
20,000			X	15.08
50,000			X	30.17

V. CONCLUSIONS

This paper discusses the design, development and testing of workflows and web services for shrimp SNPs analysis using Taverna and myGrid Technology. Local web services communicating with a BLAST application solve time-out problems and significantly reduce the working time compared to manual operation. The workflow successfully reveals genetic relationships for shrimp. The workflow can also help with analyzing other species, such as mollusc and arthropoda. The user can supply input sequences for the species to the workflow without any modification.

ACKNOWLEDGEMENTS

This research work is a collaboration between the Computer Engineering Department, the PSU Grid Center, and The Centre for Genomics and Bioinformatics Research at Prince of Songkla University. This work is partly supported by the Biogrid Project funded by the Thai National Grid Center and Thailand’s Software Industry Promotion Agency. The authors are grateful to Prof. Carole Goble, Prof. Andy Brass and Dr. Katy Wolstencroft at the School of Computer Science at the University of Manchester, and myGrid team members for Taverna training. The authors are also thankful to Dr. Andrew Davison for proof reading the paper.

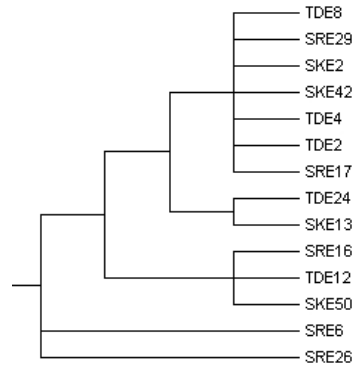


Figure 5. DNA Phylogenetic tree

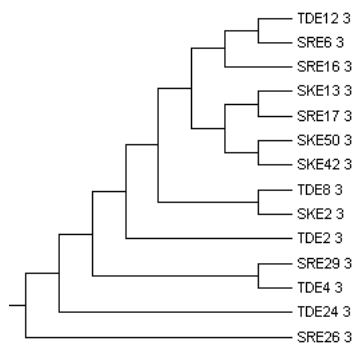


Figure 6. Amino Acid Phylogenetic tree

REFERENCES

- [1] Oinn, T., Addis, M., Ferris, J., Marvin, D., Senger, M., Greenwood, M., Carver, T., Glover, K., Pocock, M.R., Wipat, A. and Li, P., “Taverna: a tool for the composition and enactment of bioinformatics workflows,” *BIOINFORMATICS*, 2004 Vol. 20, no. 17, pp. 3045–3054, doi:10.1093/bioinformatics/bth361.
- [2] Wolstencroft, K.; Oinn, T.; Goble, C.; Ferris, J.; Wroe, C.; Lord, P.; Glover, K.; Stevens, R., “Panoply of utilities in Taverna,” *First International Conference on e-Science and Grid Computing*, 5-8 Dec. 2005.
- [3] Hull, D., Wolstencroft, K., Stevens, R., Goble, C., Pocock, M.R., Peter Li, P. and Oinn, T., “Taverna: a tool for building and running workflows of services,” *Nucleic Acids Research*, 2006, Vol. 34, Web Server issue W729–W732.
- [4] Turi, D., Missier, P., Goble, C., De Roure, D., Oinn, T., “Taverna Workflows: Syntax and Semantics, e-Science and Grid Computing,” *IEEE International Conference on 10-13 Dec. 2007*, pp. 441 – 448.
- [5] Taverna project website. 2008. <http://tavern.sourceforge.net>.
- [6] SNP Home. 2008. <http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp>.
- [7] Wanna, W., Rolland, J.L., Bonhomme, F. and Phongdara, A., “Population genetic structure of *Penaeus merguensis* in Thailand based on nuclear DNA variation,” *J. Exp. Mar. Bio, 2004, Ecol. 311*, pp. 63-78.
- [8] Phylogenetic trees. 2008. <http://cnx.org/content/m11052/latest/>.
- [9] EBI Web Services for EMBOSS-4.1 application. 2008. <http://www.ebi.ac.uk/soaplab/emboss4/>
- [10] Soaplab. 2008. <http://www.ebi.ac.uk/soaplab>.
- [11] BLAST: Basic Local Alignment and Search Tool document. 2008. <http://www.ncbi.nlm.nih.gov/blast/>.
- [12] Web API for Bioinformatics. DDBJ. 2008. <http://xml.nig.ac.jp/index.html>.

- [13] myGrid. 2008. <http://www.mygrid.org.uk/>.
- [14] Biological web services. 2008. <http://www.mygrid.org.uk/wik/myGrid/BiologicalWebServices>.
- [15] Using ClustalX. 2008. http://www.bioinformatics.ubc.ca/resources/tutorials/using_clustalx.
- [16] GeneDoc. 2008. <http://www.nrbcs.org/gfx/genedoc/>.
- [17] PSA Structure Prediction. 2008. <http://bmerc-www.bu.edu/psa/>.
- [18] BioEdit Sequence Alignment Editor for Windows 95 /98 /NT /XP. 2008. <http://www.mbio.ncsu.edu/BioEdit/BioEdit.html>
- [19] The PSIPRED protein structure prediction server. 2008. <http://bioinf.cs.ucl.ac.uk/psipred/>.
- [20] BLAST. XML Center of DDBJ. 2008. <http://www.xml.nig.ac.jp/wabi/Method?serviceName=Blast&mode=methodList&lang=en>
- [21] EMBOSS explorer. 2008. <http://pro.genomics.purdue.edu/cgi-bin/emboss>.
- [22] FASTA format description. 2008. <http://www.ncbi.nlm.nih.gov/blast/fasta.shtml>.
- [23] Hanuman. Grid Center of Prince of Songkla University website. 2008. <http://hanuman.psu.ac.th/>.
- [24] JWS. Web Services – Axis. 2008 <http://ws.apache.org/axis>.
- [25] PHYLIP Home Page. 2008. <http://evolution.genetics.washington.edu/phylip.html>
- [26] Rocks Clusters. 2008. <http://www.Rocksclusters.org/>.
- [27] Apache Tomcat. 2008. <http://tomcat.apache.org/download-55.cgi>.
- [28] AJAX Command Definition. 2008. <http://emboss.sourceforge.net/developers/acd/>.

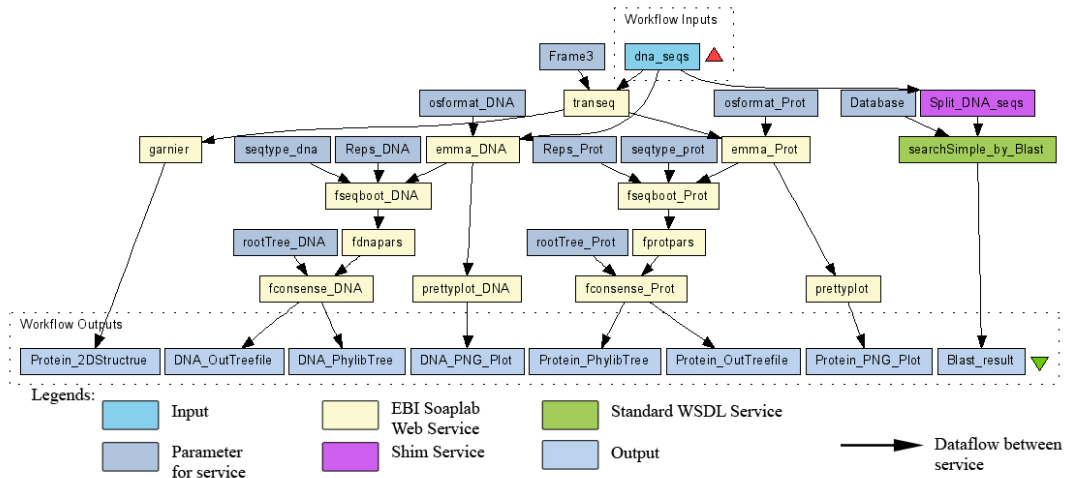
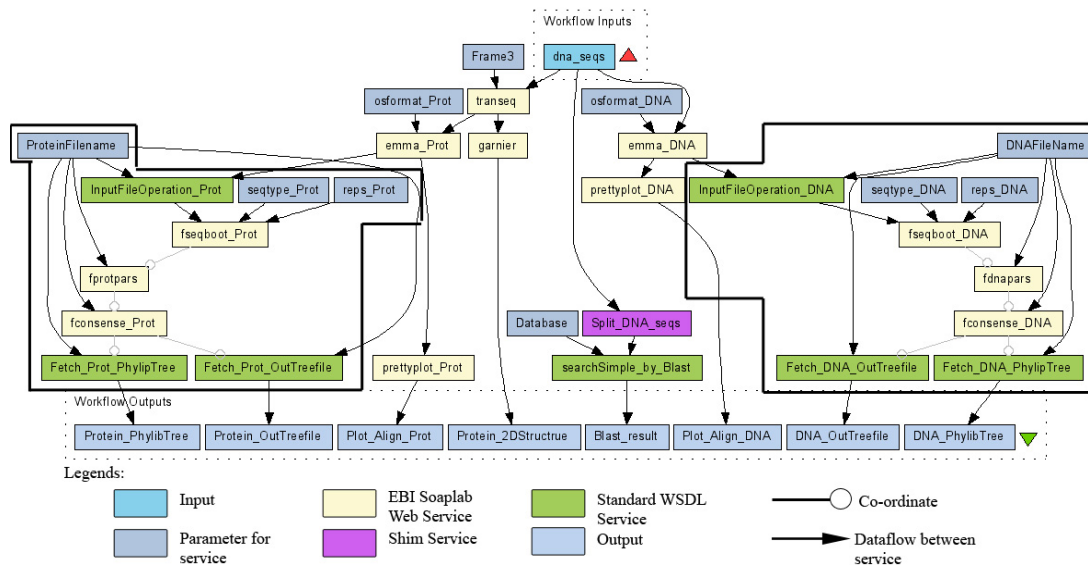


Figure 2. The first workflow of the flowchart



The new workflow uses local web services is in the black frame

ภาคผนวก ข

บทความวิชาการระดับนานาชาติชิ้นที่ 2

เรื่อง

Taverna Workflow for Validating Biomart Services

งานประชุม

2009 International Conference on Computer Design and Applications

(ICCDA 2009), Singapore, 15th - 17th May 2009.

PP. 459-463. ISBN: 978-1-4244-4213-3

Taverna Workflow for Validating BioMart Services

Kasikrit Damkliang¹ and
Pichaya Tandayya²

Department of Computer Engineering
Faculty of Engineering
Prince of Songkla University
Hat Yai, Songkhla, Thailand
s5010120123@psu.ac.th¹,
pichaya@coe.psu.ac.th²

Ekawat Pasomsub³
and Wasun Chantratita⁴

Department of Clinical Pathology
Faculty of Medicine Ramathibodi
Hospital, Mahidol University
Bangkok, Thailand
g4737429@student.mahidol.ac.th³,
rawct@mahidol.ac.th⁴

Surakameth Mahasirimongkol
Center for International Cooperation
Department of Medical Sciences
Ministry of Public Health
Nonthaburi, Thailand
surakameth.m@dmisc.mail.go.th

Abstract—This paper present a novel Taverna workflow for validating BioMart Services. It saves a significant amount of time by avoiding incorrect workflows caused by out-of-date services. It also reduces debugging time by finding out-of-date processor query fields (filters) and output fields (attributes). In addition, it provides associated information required for future validation.

Keywords—*workflow; BioMart; Taverna; Web Services; validation*

VI. INTRODUCTION

Workflow is a data flow methodology where output of one component forms the input of others [1][6]. If any service in the workflow produces erroneous results, then the entire data flow in that section will produce incorrect results.

To decide whether a workflow gives correct results or not, we normally need to run the workflow and analyse its results. If the workflow is complex, with a lot of data to analyse, then it will take a long time to produce results. The results cannot be guaranteed until they are verified with other research work or data [2]. Therefore, validating the workflow before running it would help avoid unnecessary and non-productive processing, and reduce the working time.

Taverna version 1.7.0 (and later) has a workflow existence plug-in, which can check responsiveness of web service endpoints [3][4]. However, it cannot detect upgrade changes to web services, which produce empty result lists, even though the endpoints are correct.

This problem also impacts the search for processor query fields (filters) and processor output fields (attributes) which may simply be out-of-date [5]. The more complex the workflow, the more time will be lost in running and debugging the workflow.

The test case described here is a complex pharmacogenomic research workflow, requiring BioMart Web services to gather information from databases. The BioMart system is a flexible data warehouse aimed at complex interlinked biological data sets. Taverna's BioMart query integration provides full search and retrieval functionality over these data sources [4][7].

This paper presents a workflow method based on the Taverna environment which validates the BioMart services before they are utilized in the processor query fields and

processor output fields. This detects out-of-date services, filters and attributes, greatly reducing debugging times.

VII. BACKGROUND

Many useful bioinformatic databases are accessible via the Internet, but finding the most suitable for a task is difficult since there are a few standards for representing and sharing data. myGrid e-Science provides an analysis method depending on a workflow system which can convert data to comply with web service interfaces and direct the flow of data between resources [1][6][8]. The myGrid e-Science Taverna tool is an open source workflow composer for orchestrating bioinformatic web services and existing applications into workflows. In essence, Taverna helps compose local and distributed resources into a single logical system, by employing Web Service technology [4]

Taverna can access over 30,000 services or processors in the bioinformatics domain, and there are plans to include other fields such as astronomy, chemoinformatics, and health informatics [9].

Taverna can add service endpoints, such as Web Service Definition Language (WSDL) documents [10], for improved portability, and can access other types of resources, such as, local Java services, BioMoby, beanshell scripts, and BioMart services [4][7].

BioMart [7][11] is a widely used query-oriented data management system jointly developed by the Ontario Institute for Cancer Research (OICR) and the European Bioinformatics Institute (EBI). It is open source and freely distributed without restrictions. The system accepts various types of data and also provides a kind of 'data mining.' Service configuration is achieved through a graphical user interface based on applications using web services, a dedicated API. A BioMart service configuration is saved into a Simple Conceptual Unified Flow Language (Scufl) [4] document in the workflow. BioMart is available to Taverna as a plug-in, providing good ways to compose workflows together.

Taverna version 1.7.0, and later, has a 'Taverna 2 preview' plug-in for performing workflow health checks. The plug-in examines whether the workflow can be translated from the old style Scufl model into the T2 style, and whether the web service endpoints are responding [4]. It

cannot detect upgrade changes for BioMart processor filters and attributes, or examine services inside nested workflows. As a consequence, 'Taverna 2 preview' is somewhat lacking for validating services and reducing debugging time.

The following running and debugging steps are described in order to show how difficult it is for a user to check and debug a workflow by themselves.

When a user knows a workflow has produced incorrect results, he will check the data flow by verifying the input and output of each service and processor, which should highlight the erroneous resource. If not, then the data type and syntax of the resources must be examined. Even if these seem alright, the workflow may still produce incorrect results.

In that case, the problem may actually lie within the processor. Only then will the out-of-date processor filters and processor attributes be found.

Clearly, the main problem is when the BioMart services are upgraded, causing workflows to produce incorrect data flow.

The out-of-date filters or attributes can be found manually by looking for suspicious resources on the configuration screen for the BioMart processor node. If certain filters or attributes are absent then, it means they are out-of-date. However, manual validation can take many hours or days.

VIII. BIOMART WORKFLOW VALIDATION

A. Design

A workflow composed in the Taverna environment is saved into a Scufi document, a standard format from myGrid. The Scufi employs an XML (Extensible Markup Language) [12] so documents can be used with the XML Path Language (XPath) [13]. XPath is employed to extract data elements from the Scufi of a BioMart processor.

The BioMart processor is configured as a 'data mining' resource, which consists of a processor, or service name, and a dataset. The dataset is made of information extracted from the configuration databases; for example, the 'Homo sapiens SNPs' dataset takes data from the dbSNP127, ENSEMBL, TSC1 and HGVbase15 databases [7][8].

The BioMart dataset configuration is presented with the query fields (filters) and output fields (attributes) matching the user's requirements [5]. The BioMart service structure is shown in Figure 1.

A BioMart Service		
A Processor/Service Name		
A Dataset	Query fields (Filters)	Output field (Attributes)

Figure 3. The BioMart service structure in a Taverna workflow

An example Scufi workflow implementation of a BioMart processor is shown below:

```
<s:processor name="hsapiens_gene_GO">
  <biomart:query>
    <biomart:Dataset name="hsapiens_gene_ensembl">
      <biomart:Attribute name="go" />
      <biomart:Attribute name="evidence_code" />
      <biomart:Attribute name="go_description" />
    </biomart:query>
  </s:processor>
```

```
<biomart:Attribute name="ensembl_gene_id" />
<biomart:Filter name="ensembl_gene_id" value="" list="true"/>
</biomart:Dataset>
</biomart:Query>
</s:processor>
```

When Taverna opens a workflow, the responsiveness of the endpoints will be checked. If any problematic endpoints are found then Taverna will issue a warning message, and not run the workflow [4]. Since Taverna handles the endpoints, our work concentrates on validating the dataset filters and attributes in the processors. A data flow context diagram for the BioMart validating workflow is shown in Figure 2.

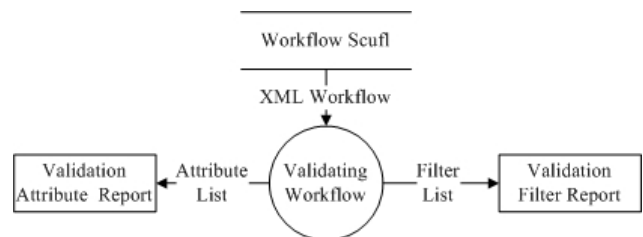


Figure 4. A data flow context diagram for BioMart debugging workflow

The Scufi workflow document(s) fed into the *Validating Workflow* process will produce a filter and attribute validation report.

The processor name used in Figure 3 is extracted from the Scufi document, and is vital for obtaining the dataset name. Subsequently, the Scufi document, processor names and dataset names are fed into a loop of validating processes, as shown in Figure 4.

Fig. 4 describes how out-of-date filters and attributes can be detected by checking their availability with the central BioMart dataset registry on the Internet. If any attribute or filter has been updated, the system collects its description, and details of each processor. This validation report will be utilized to cross-check for BioMart upgrades.

B. Implementation

The nested workflow in Figure 4 acts as an iteration module for the main task. It is utilized to check all the BioMart processor filters and attributes in Scufi documents against the central BioMart registry, located on the BioMart web services server (<http://www.BioMart.org/BioMart/martservice>).

The metadata about service availability is obtained via a GET request using a 'Get web page from URL' Java local service. A number of metadata queries are available, which can be formulated by appending parameters to the end of the URL [11]. The first parameter is appended using the ? symbol, and subsequent parameters use the & symbol. The format is:

```
../martservice?type=<filters|attributes>&virtualschema=de
fault &dataset=<datasetname>
```


For example, the command for retrieving the attribute registry of the 'hsapiens_gene_ensembl' BioMart is:

http://www.BioMart.org/BioMart/martservice?type=attribute&virtualschema=default&dataset=hsapiens_gene_ensembl

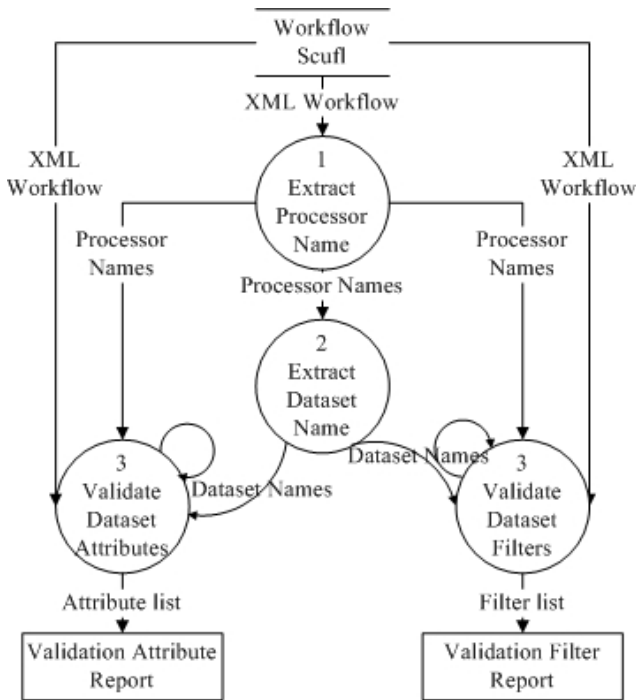


Figure 5. The second level data flow diagram of our validating process

XPath expressions are utilized to extract data elements from the Scuff workflow using the 'Xpath From Text' Java local service, as shown in Table I. The validating processes are coded in Java using the Beanshell local service [3].

Fig. 5 shows the composed workflow for validating the filters and attributes of large BioMart services. The 'Check_Filter' and 'Check_Attribute' rectangular frames are nested looping workflows.

Taverna manages the data flow of output ports using the cross product ('all against all') relationship. In our work, the data flow management for integrating the nested workflows into the upper processes is carried out by the Taverna 'Configure Iterators' [3] as shown in Figure 6. The processor name and dataset name are used in a dot product ('one against one') relationship, whereas the workflow Scuff employs cross product. The relationship output of the data ports is shown in Table II, showing that a Scuff has many processors, and each processor has only one dataset.

IX. EXPERIMENTAL RESULTS

Our first test utilized an out-of-date pharmacogenomic workflow that no longer produces a correct data flow, and we focused on reducing the validating time and avoiding

wasteful executions of workflows that produce no results. The simple testing workflow for Gene Ontology [14] is shown in Figure 7, utilizing a BioMart processor 'hsapiens_gene_ensembl' labelled as 'hsapiens_gene_GO.' The service consists of one filter and four attributes. We ran the workflow using two Ensembl IDs. The response arrived in less than two seconds, but was empty.

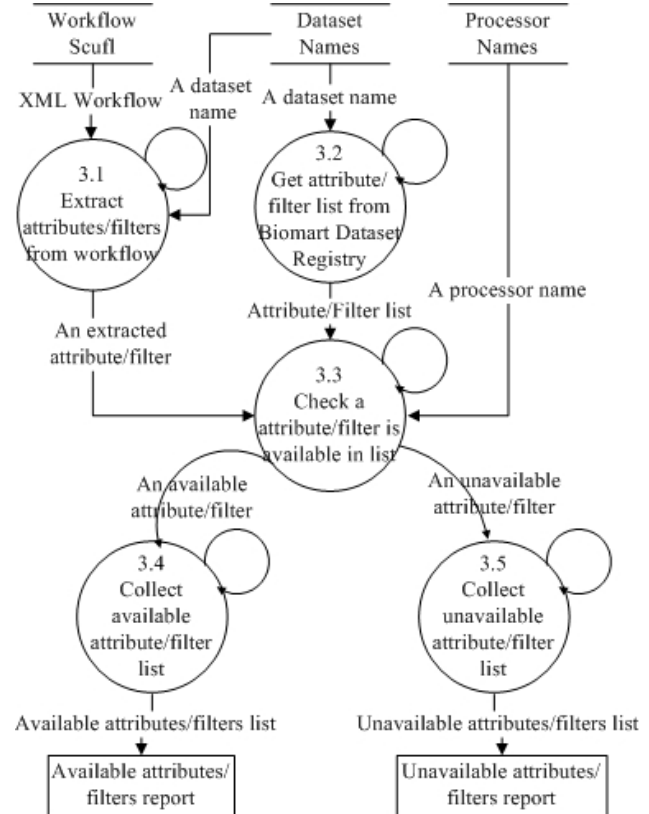


Figure 6. The third level data flow diagram of our validating workflow for dataset filters/attributes

TABLE VI. XPATH EXPRESSION IN THE WORKFLOW

Extracting	Xpath expression
Processor names	<code>//*[local-name()='Dataset' and (namespace-uri()='http://org.embl.ebi.escience/xScuff-Biomart/0.1alpha')/ancestor::s:processor/@name</code>
BioMart processor names	<code>//*[local-name()='processor'] and (@name=<A BioMart processor name>)]//*[local-name()='biomart']/parent::s:processor/@name</code>
BioMart dataset names	<code>//*[local-name()='processor'] and (@name=<A BioMart processor name>)]//*[local-name()='biomart']//*[local-name()='MartQuery']//*[local-name()='Query']//*[local-name()='Dataset']/@name</code>
BioMart dataset filters or dataset attributes	<code>//*[local-name()='processor'] and (@name=<A BioMart processor name>)]//*[local-name()='biomart']//*[local-name()='MartQuery']//*[local-name()='Query']//*[local-name()='Dataset']//*[local-name()='Filter' 'Attribute']/@name</code>

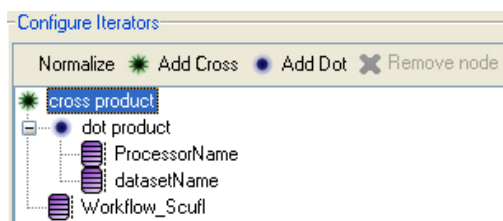


Figure 6. Configure Iterators in Taverna

TABLE VII. THE RESULTS OF CONFIGURE ITERATORS

Cross product	Dot product	Dot product
Workflow_Scuff	1 st ProcessorName	1 st datasetName
Workflow_Scuff	2 nd ProcessorName	2 nd datasetName

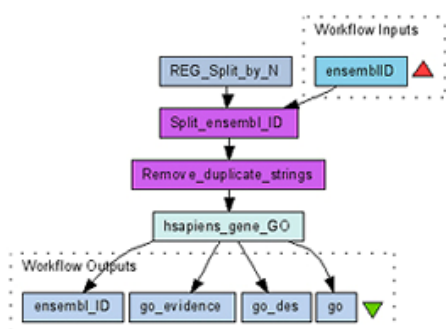


Figure 7. An example of a testing workflow

Our validating workflow helps to identify the out-of-date filters and attributes, as shown below:

Processor Name: *hsapiens_gene_GO*, Attribute: *go* is out-of-date
 Processor Name: *hsapiens_gene_GO*, Attribute: *evidence_code* is out-of-date
 Processor Name: *hsapiens_gene_GO*, Attribute: *go_description* is out-of-date

The validating workflow also collects the up-to-date filters and attributes for later use:

Processor Name: *hsapiens_gene_GO*, Filter: *ensembl_gene_id* is up-to-date
 Detail: Filter to include genes with supplied list of Ensembl Gene IDs
 Processor Name: *hsapiens_gene_GO*, Attribute: *ensembl_gene_id* is up-to-date
 Detail: Ensembl Stable ID of the Gene, feature_page

The description of the '*ensembl_gene_id*' attribute, is the Ensembl Stable ID for the Gene, and is available on the feature page of the BioMart service configuration screen.

The second major test of our validating workflow is a pharmacogenomic workflow, available at <http://www.myexperiment.org/workflows/610>.

The workflow testing results are shown in Table III. It consists of 52 services, including 17 that link to BioMart services, 70 corresponding data links, and three data querying BioMart services with three filters and 22 attributes.

The workflow inputs are 500 to 1,000 Single Nucleotide Polymorphisms (SNPs) are the smallest genetic statistic P values according to the Genome Wide Association Study (GWAS) [2] of 550,000 SNPs.

It took 14 hours to run the workflow with 1,000 SNP IDs. Every service works successfully except the Gene Ontology part which returns an empty list.

As shown in Table III, it took 26.5 seconds to validate the workflow, and identify the out-of-date filters and attributes at the Gene Ontology finding service. This shows that validating a workflow before executing it can significantly reduce work when there are out-of-date services. Debugging time for finding the cause of incomplete results is also reduced. After the validation, we can configure our workflow so that all the units are up-to-date, as shown at <http://www.myexperiment.org/workflows/612>.

TABLE VIII. THE PHARMACOGENOMIC WORKFLOW TESTING RESULT

Type	Property	Remark
All Service	52	-
All Data link	152	-
Associated BioMart service	17	3 BioMart query services
Associated BioMart data link	70	-
BioMart Filter	3	3 BioMart query services
BioMart Attribute	22	3 BioMart query services
Test input	1000	SNP IDs
Running time	14.00.50	Hours
Running results	All completed	Except Gene Ontology
Validating time	26.5	Seconds

X. RESULT DISCUSSION AND DEBUGGING SOLUTION

The experimental results show that minor workflow re-design is often required: out-of-date services must be removed and new services added.

Even when there is only one out-of-date filter, the workflow will fail to produce correct data flows because query fields must pass all the filter conditions before retrieving information. The user must fix the offending processor node before the workflow will function.

Even a processor with an out-of-date attribute that is not an input for downstream processes may still cause the processor to produce incorrect data.

Out-of-date filters and attributes due to the upgrades are not shown on the BioMart configuration screen, so the user cannot deselect them. Advanced users may manually edit the Scuff document but this tends to be tedious, confusing and error-prone.

Our advice is to add a new service with the same processor name that occurred in the validation report. This is done by reselecting the service from Taverna's service available panel, and then attaching the processor and its data links to the other processors. When the added processor does not have the necessary filters or attributes, then replacement services in the same domain should be utilized if possible.

XI. CONCLUSIONS

We have developed a Taverna workflow for validating the processor query fields and processor output fields in BioMart services. This validating workflow can help improve workflow processes without requiring the execution of the workflow. It can save significant amounts of computing time when running incorrect workflows, and

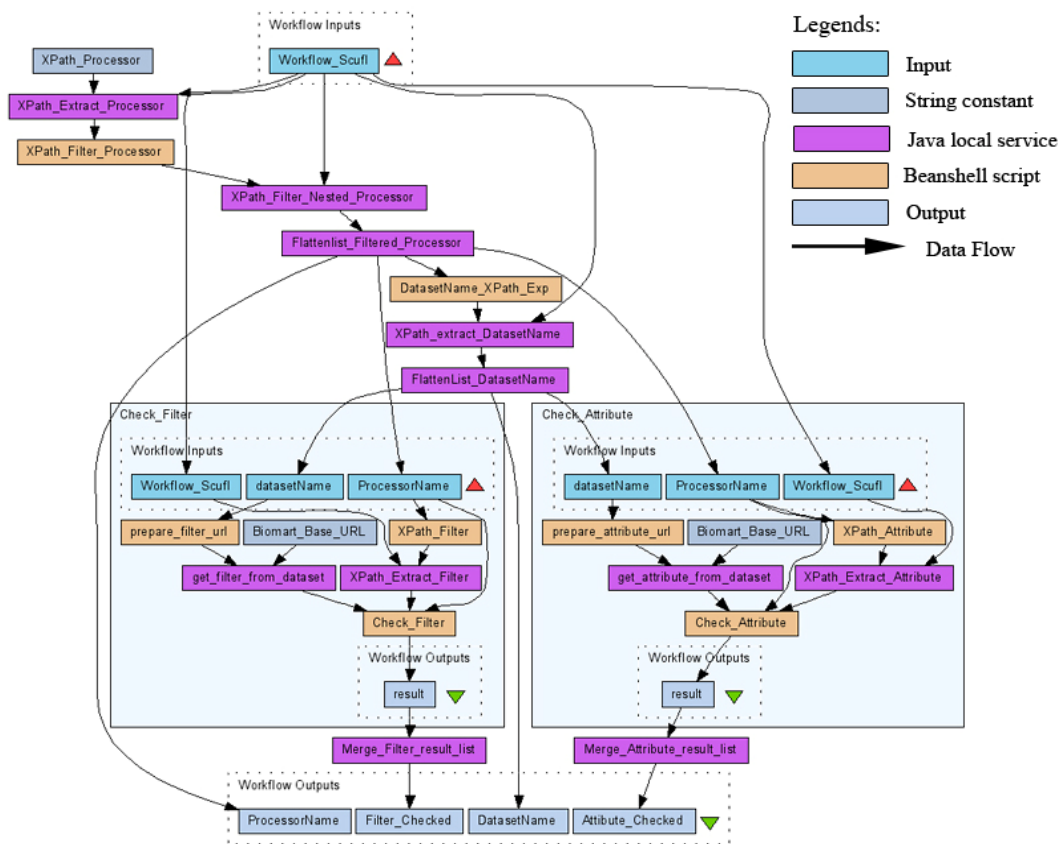
reduces the debugging time. It also provides the associated information useful for future validations.

ACKNOWLEDGMENTS

This research work is a collaboration between the Computer Engineering Department, the PSU Grid Center at Prince of Songkla University, and the Thailand Center of Excellence for Life Sciences (TCELS) at Mahidol University. This work is partly supported by the Biogrid Project funded by the Thai National Grid Center (TNGC), Thailand's Software Industry Promotion Agency (SIPA) and TCELS. The authors are grateful to Prof. Carole Goble, Prof. Andy Brass and Dr. Katy Wolstencroft at the School of Computer Science at the University of Manchester, and myGrid team members for Taverna training. The authors are also thankful for Dr. Andrew Davison for proof reading the paper.

REFERENCES

- [1] Hull, D., Wolstencroft, K., Stevens, R., Goble, C., Pocock, M.R., Peter Li, P. and Oinn, T., "Taverna: a tool for building and running workflows of services," Nucleic Acids Research, 2006, Vol. 34, Web Server issue W729–W732.
- [2] Wellcome Trust Case Control Consortium. "Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls," Nature, 2007;447, pp. 661-678.
- [3] Taverna project website. 2009. <http://taverna.sourceforge.net/>.
- [4] myGrid. 2009. <http://www.mygrid.org.uk/>.
- [5] Rice, P.M.; Bleasby, A.J.; Haider, S.A.; Ison, J.C.; Meglinchey, S.; Uludag, M., "EMBRACE: Bioinformatics Data and Analysis Tool Services for e-Science," e-Science and Grid Computing, 2006, e-Science '06. Second IEEE International Conference on Dec. 2006, pp. 146 – 146, Digital Object Identifier 10.1109/E-SCIENCE.2006.261079.
- [6] Turi, D., Missier, P., Goble, C., De Roure, D., Oinn, T., "Taverna Workflows: Syntax and Semantics, e-Science and Grid Computing," IEEE International Conference on 10-13 Dec. 2007, pp. 441 – 448.
- [7] BioMart. 2009. <http://www.BioMart.org/>.
- [8] Belhajjame, K.; Wolstencroft, K.; Corcho, O.; Oinn, T.; Tanoh, F.; William, A.; Goble, C., "Metadata Management in the Taverna Workflow System," Cluster Computing and the Grid, 2008. CCGRID '08. 8th IEEE International Symposium on 19-22 May 2008, pp. 651 – 656, Digital Object Identifier 10.1109/CCGRID.2008.17.
- [9] Biological web services. 2009. <http://www.mygrid.org.uk/wiki/Mygrid/BiologicalWebServices>.
- [10] Web Services Description Language (WSDL). 2008. <http://www.w3.org/TR/wSDL>.
- [11] Smedley, D., Haider, S., Ballester, B., Holland, R., London, D., Thorisson, G. and Kasprzyk, A., "BioMart - biological queries made easy," BMC Genomics 2009, 10:22doi:10.1186/1471-2164-10-22.
- [12] Extensible Markup Language (XML). 2009. <http://www.w3.org/XML/>.
- [13] XML Path Language (XPath). 2009. <http://www.w3.org/TR/xpath>.
- [14] The Gene Ontology. 2009. <http://www.geneontology.org/>



The workflow for validating the BioMart filters and BioMart attributes

ภาคผนวก ก

สื่อสาธิตวิธีการทำงานของเวิร์คโฟลว์สำหรับตรวจสอบบริการไปรษณีย์

(คู่มือที่แนบมากับวิทยานิพนธ์)