



ใบรับรองวิทยานิพนธ์

บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์

วิศวกรรมศาสตรมหาบัณฑิต (วิศวกรรมคอมพิวเตอร์)

ปริญญา

วิศวกรรมคอมพิวเตอร์

วิศวกรรมคอมพิวเตอร์

สาขา

ภาควิชา

เรื่อง ขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพโดยพิจารณาความถี่ส่วนย่อยของลำดับ

A Sub-Segment Frequency-Based Multiple Alignment Algorithm for Biological Sequences

นามผู้วิจัย นายคมสัน จันมา

ได้พิจารณาเห็นชอบโดย

ประธานกรรมการ

(ผู้ช่วยศาสตราจารย์พันธุ์ปิติ เปี่ยมสง่า, D.Sc.)

กรรมการ

(ผู้ช่วยศาสตราจารย์พีรวัฒน์ วัฒนพงษ์, Ph.D.)

กรรมการ

(ผู้ช่วยศาสตราจารย์เข็มชาติ วิภาตะวนิช, Ph.D.)

หัวหน้าภาควิชา

(ผู้ช่วยศาสตราจารย์เข็มชาติ วิภาตะวนิช, Ph.D.)

บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์รับรองแล้ว

(รองศาสตราจารย์วินัย อัจจงหาญ, M.A.)

คณบดีบัณฑิตวิทยาลัย

วันที่ เดือน พ.ศ.

วิทยานิพนธ์

เรื่อง

ขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพโดยพิจารณาความถี่ส่วนย่อยของลำดับ

A Sub-Segment Frequency-Based Multiple Alignment Algorithm for Biological Sequences

โดย

นายคมสัน จันมา

เสนอ

บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์

เพื่อความสมบูรณ์แห่งปริญญาวิทยาศาสตรมหาบัณฑิต (วิศวกรรมคอมพิวเตอร์)

พ.ศ. 2550

คมสัน จันมา 2550: ขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพโดยพิจารณา
ความถี่ส่วนย่อยของลำดับ ปรินญาวิศวกรรมศาสตรมหาบัณฑิต (วิศวกรรม
คอมพิวเตอร์) สาขาวิศวกรรมคอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์ ธรรมศาสตร์
มหาวิทยาลัยที่ปรึกษา: ผู้ช่วยศาสตราจารย์พันธุ์ปิติ เปี่ยมสง่า, D.Sc. 100 หน้า

การเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพเป็นงานที่จำเป็นสำหรับนักชีววิทยา ซึ่งใน
ปัจจุบันจำนวนของข้อมูลทางด้านชีวภาพได้เพิ่มขึ้นอย่างรวดเร็ว ดังนั้นเครื่องมือสำหรับการเทียบ
เรียงกลุ่มลำดับข้อมูลทางชีวภาพจึงจำเป็นต้องมีความถูกต้องในระดับที่นักชีววิทยาพอใจ และใช้
เวลาในการประมวลผลน้อย งานวิจัยนี้จึงนำเสนอวิธีการเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพที่
ใช้เวลาในการทำงานน้อยมาก โดยที่ความถูกต้องเทียบเคียงได้กับเครื่องมือที่ใช้ในปัจจุบัน โดยใช้
วิธี n-gram ตัดข้อมูลเป็นส่วนย่อย และพิจารณาจากความถี่ของส่วนย่อยของแต่ละคู่ลำดับข้อมูล
เพื่อใช้เป็นแนวทางในการเทียบเรียง วัดความถูกต้อง และเวลาการประมวลผลโดยใช้ชุดข้อมูล
ทดสอบ BALiBASE และ PREFAB ผลลัพธ์ที่ได้จากการเทียบเรียงชุดทดสอบ PREFAB สามารถ
เทียบเรียงได้เร็วกว่า ClustalW 6.24 เท่า โดยที่ความถูกต้องอยู่ในระดับเดียวกัน

Comsan Chanma 2007: A Sub-Segment Frequency-Based Multiple Alignment Algorithm for Biological Sequences. Master of Engineering (Computer Engineering), Major Field: Computer Engineering, Department of Computer Engineering. Thesis Advisor: Assistant Professor Panpiti Piamsa-nga, D.Sc. 100 pages.

Multiple sequence alignment has become a necessary task of biologists for sequence analysis. Because of fast increasing of biology data, biologists not only need accurate tools but fast tools are also required for multiple sequence alignment of biology data. In this paper, we propose a sub-segment's frequency-based processing technique to improve speed of alignment process. This technique is pre-process a sub-segment's frequency matrix between all pair of sequences with n-gram technique. This provides us with a distance matrix that can be used to guide the progressive alignment. Accuracy and computing time are measured using BALiBASE and PREFAB benchmark. The CPU time is drastically reduced as compared with ClustalW. On the PREFAB benchmark alignment databases, our result is 6.24 times faster than ClustalW with comparable accuracy.

Student's signature

Thesis Advisor's signature

____ / ____ / ____

กิตติกรรมประกาศ

ข้าพเจ้าขอขอบพระคุณ ผู้ช่วยศาสตราจารย์พันธุ์ปิติ เปี่ยมสง่า ประธานกรรมการที่ปรึกษา
ที่ให้คำแนะนำและให้ความช่วยเหลือในการทำวิทยานิพนธ์ฉบับนี้เป็นอย่างสูง ขอขอบพระคุณ
ผู้ช่วยศาสตราจารย์พีรวัฒน์ วัฒนพงศ์ กรรมการสาขาวิชาเอก ผู้ช่วยศาสตราจารย์ขณะทัต วิชา
ตะวนิช กรรมการสาขาวิชารอง และผู้ที่เกี่ยวข้องกับวิทยานิพนธ์ฉบับนี้ที่ให้คำแนะนำเกี่ยวกับการ
นำเสนองานวิจัย และช่วยตรวจสอบแก้ไขข้อบกพร่องต่าง ๆ ของวิทยานิพนธ์ฉบับนี้

ขอขอบคุณคุณคุณจิตติมนต์ เขียนดวงจันทร์ และคุณวีรวุฒิ คงบุญเกียรติ ที่ให้แนวคิดและ
แนวทางในการทำงานวิจัย และขอบคุณสมาชิกในห้องปฏิบัติการวิจัยการวิเคราะห์และค้นพบ
เนื้อหาประสมที่เอื้อเพื่อสถานที่ในการทำวิทยานิพนธ์ และให้กำลังใจในการทำวิทยานิพนธ์เสมอมา

สุดท้ายนี้ขอขอบคุณคุณพ่อ คุณแม่ และทุกคนในครอบครัวที่คอยเป็นกำลังใจที่ดีที่สุด
ขอบคุณรุ่นพี่ รุ่นน้อง และเพื่อนๆทั้งที่วิศวกรรมศาสตร์ มหาวิทยาลัยเกษตรศาสตร์ และที่บริษัท ที.
เอ็น. อินฟอร์เมชั่น ซิสเต็มส์ จำกัด ที่ให้ความช่วยเหลือ และเป็นกำลังใจผลักดันให้วิทยานิพนธ์
ฉบับนี้เสร็จสมบูรณ์ได้

คมสัน จันมา

มีนาคม 2550

สารบัญ

	หน้า
สารบัญ	(1)
สารบัญตาราง	(2)
สารบัญภาพ	(3)
คำนำ	1
วัตถุประสงค์	2
การตรวจเอกสาร	3
ความรู้เบื้องต้นเกี่ยวกับข้อมูลชีวภาพ	3
รูปแบบการจัดเก็บข้อมูลชีวภาพ	8
ลักษณะการเปลี่ยนแปลงของข้อมูลชีวภาพ	8
การเทียบเรียงข้อมูล	9
ประโยชน์ของการเทียบเรียงกลุ่มข้อมูลชีวภาพ	12
ตารางค่าคะแนน	13
เทคนิคที่ใช้ในการเทียบเรียงกลุ่มข้อมูล	14
โครงสร้างการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW	18
ปัญหาของการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพ	20
อุปกรณ์และวิธีการ	22
อุปกรณ์	22
วิธีการ	23
ผลและวิจารณ์	46
สรุปและข้อเสนอแนะ	60
สรุป	60
ข้อเสนอแนะ	61
เอกสารและสิ่งอ้างอิง	62
ภาคผนวก	64
ประวัติการศึกษา และการทำงาน	100

สารบัญตาราง

ตารางที่		หน้า
1	ตารางค่าคะแนน BLOSUM50	16
2	ตารางค่าตัวเลขสำหรับแทนค่าอักขระในส่วนย่อยของลำดับ	24
3	ตารางค่าความแตกต่างของลำดับ a, b และ c ที่เข้าเทียบเรียง	26
4	ชุดอ้างอิง 1 เป็นชุดที่มีค่าความห่างเท่ากันและหลายระดับส่วนสงวน (ขนาดสั้น)	30
5	ชุดอ้างอิง 1 เป็นชุดที่มีค่าความห่างเท่ากันและหลายระดับส่วนสงวน(ขนาดกลาง)	30
6	ชุดอ้างอิง 1 เป็นชุดที่มีค่าความห่างเท่ากันและหลายระดับส่วนสงวน (ขนาดยาว)	31
7	ชุดอ้างอิง 2 เป็นการเทียบเรียงกับสายอักขระที่มีการเปลี่ยนแปลงสูง	31
8	ชุดอ้างอิง 3 เป็นกลุ่มย่อยที่น้อยกว่า 25% ของส่วนที่เหลือระหว่างกลุ่ม	31
9	ชุดอ้างอิง 4 และชุดอ้างอิง 5	32
10	ชุดอ้างอิง 1	32
11	ชุดอ้างอิง 2-5	33
12	กลุ่มข้อมูลที่ 1-7	34
13	กลุ่มข้อมูลที่ 8-14	35
14	กลุ่มข้อมูลที่ 15 – 21	36
15	กลุ่มข้อมูลที่ 22 - 28	37
16	กลุ่มข้อมูลที่ 29 - 35	38
17	กลุ่มข้อมูลที่ 36 - 42	39
18	กลุ่มข้อมูลที่ 43 - 49	40
19	กลุ่มข้อมูลที่ 50 - 56	41
20	กลุ่มข้อมูลที่ 57 - 63	42
21	กลุ่มข้อมูลที่ 64 - 68	43

สารบัญภาพ

ภาพที่		หน้า
1	โครงสร้างโมเลกุลของเบสอันเป็นองค์ประกอบของนิวคลีโอไทด์	4
2	โครงสร้างของดีเอ็นเอ เป็นแบบเกลียวคู่ขนาน (double helix)	5
3	โครงสร้างทั่วไปทางเคมีของกรดอะมิโน	6
4	โครงสร้างโมเลกุล และชื่อของกรดอะมิโนทั้ง 20 ชนิด	7
5	ตัวอย่างการจัดเก็บข้อมูลชีวภาพแบบ FASTA	8
6	รูปแบบการเปลี่ยนแปลงของเบสในดีเอ็นเอ โดย 1) คือลำดับเบสของดีเอ็นเอสายปกติ 2) คือการเปลี่ยนจาก C ไป T เป็นการเปลี่ยนของคนละคู่เบสกัน(Transition) 3) คือการเปลี่ยนจาก G ไป C เป็นการเปลี่ยนกับคู่คอมพลิเมนต์กัน (Transversion) 4) คือการขาดไปของช่วงลำดับเบส ACCTA 5) คือการเพิ่มเข้ามาของช่วงลำดับเบส AAAGC 6) คือการอินเวอร์สชั่น จาก 5-GCAAAC-3 เปลี่ยนเป็น 5-GTTTGC-3	9
7	ตัวอย่างการเทียบเรียงลำดับ 2 ลำดับ	10
8	ตัวอย่างแผนภูมิต้นไม้แสดงความสัมพันธ์ทางวิวัฒนาการที่วิเคราะห์ได้จากการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพ	12
9	ตัวอย่างผลของการเทียบเรียงกลุ่มลำดับข้อมูลเพื่อระบุส่วนจำเพาะของกลุ่มลำดับข้อมูล โดยส่วนที่มีเครื่องหมาย * คือส่วนจำเพาะของกลุ่มลำดับข้อมูล	13
10	การคำนวณค่าในตารางการคำนวณ	15
11	การเทียบเรียงแบบโกลบอลโดยใช้การโปรแกรมพลวัต และผลลัพธ์ที่ได้	17
12	การเทียบเรียงแบบโลคอลโดยใช้การโปรแกรมพลวัต และผลลัพธ์ที่ได้	17
13	ขั้นตอนในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW	19
14	การเทียบเรียงกลุ่มลำดับข้อมูลโดยพิจารณาความถี่ส่วนย่อยของลำดับ	23
15	ขั้นตอนการสร้างตารางค่าความแตกต่างโดยพิจารณาความถี่ส่วนย่อย	28
16	ขั้นตอนการสร้างตารางค่าความแตกต่างโดยพิจารณาความถี่ส่วนย่อย	29
17	เวลาที่โปรแกรม ClustalW และ โปรแกรม SSFA ใช้ในการเทียบเรียงกลุ่มลำดับข้อมูลของชุดอ้างอิงทั้ง 5 ชุด	47
18	ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และ โปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 1	48

สารบัญภาพ (ต่อ)

ภาพที่		หน้า
32	ค่าคะแนน SP เฉลี่ยในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA ของแต่ละชุดอ้างอิง	55
33	ค่าคะแนน TC เฉลี่ยในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA ของแต่ละชุดอ้างอิง	55
34	เวลาที่โปรแกรม ClustalW และโปรแกรม SSFA ใช้ในการเทียบเรียงฐานข้อมูล PREFAB	56
35	ค่าคะแนน Q เฉลี่ยในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA	57
36	เวลาที่ใช้ในการเทียบเรียงของโปรแกรม ClustalW และโปรแกรม SSFA ในจำนวนลำดับที่ต่างกัน	58
37	เวลาที่ใช้ในขั้นตอนที่ 1 ของการเทียบเรียงของโปรแกรม ClustalW และโปรแกรม SSFA ในจำนวนลำดับที่ต่างกัน	59
ภาพผนวกที่		
1	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 1	66
2	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 2	66
3	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 3	67
4	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 4	67
5	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 5	68

สารบัญภาพ (ต่อ)

ภาพผนวกที่		หน้า
58	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 58	94
59	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 59	95
60	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 60	95
61	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 61	96
62	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 62	96
63	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 63	97
64	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 64	97
65	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 65	98
66	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 66	98
67	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 67	99
68	ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 68	99

ขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพโดยพิจารณา ความถี่ส่วนย่อยของลำดับ

A Sub-Segment Frequency-Based Multiple Alignment Algorithm for Biological Sequences

คำนำ

ในปัจจุบันการนำศาสตร์ทางด้านคอมพิวเตอร์เข้ามาช่วยงานทางด้านชีวภาพนั้นแพร่หลายมากขึ้นเนื่องจากมีการทำงานที่รวดเร็ว และถูกต้อง ซึ่งประโยชน์ที่ได้ก่อให้เกิดความก้าวหน้าและการค้นพบสิ่งใหม่ในทางด้านชีววิทยาเป็นอันมาก การเทียบเรียงกลุ่มลำดับข้อมูลเป็นงานวิจัยด้านหนึ่งที่ได้รับความสนใจ และมีความสำคัญ เนื่องจากเป็นขั้นตอนที่จำเป็นสำหรับการวิเคราะห์ลำดับข้อมูลทางด้านชีววิทยา ซึ่งมีจำนวนข้อมูลมาก และเพิ่มขึ้นอย่างรวดเร็ว ผลที่ได้จากการเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพนั้นถูกใช้ประโยชน์ในหลายๆด้าน ได้แก่ การวิเคราะห์วิวัฒนาการของสิ่งมีชีวิต, การระบุส่วนจำเพาะของลำดับข้อมูล หรือแม้กระทั่งใช้ในการทำนายโครงสร้างของลำดับโปรตีนนั้นๆ วิธีการและขั้นตอนวิธีทางการโปรแกรมคอมพิวเตอร์ต่างๆ ถูกนำมาประยุกต์ใช้ในงานวิจัยเกี่ยวกับการเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพ เพื่อพัฒนา 2 ปัจจัยหลัก คือความถูกต้องของการเทียบเรียง และความเร็วในการประมวลผล

การเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพนั้น มีความยากระดับ NP complete การเทียบเรียงนั้นเป็นการเทียบเรียงลำดับที่มีความยาวสูง จำนวนมาก เป็นผลให้การประมวลผลใช้เวลานาน โดยที่ปกติค่าความซับซ้อนของการประมวลผลคือ $O(NL)$ เมื่อ L คือ ความยาวของแต่ละสายอักขระ และ N คือ จำนวนสายอักขระ ดังนั้นแนวทางการพัฒนาในปัจจุบันจึงมุ่งเน้นให้การเทียบเรียงใช้มีความเร็วสูงขึ้น ในขณะที่ความถูกต้องอยู่ในระดับที่ยอมรับได้

วิทยานิพนธ์นี้นำเสนอการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพในแบบวิธีฮิวริสติก โดยพิจารณาความถี่ของลำดับย่อยร่วมของลำดับข้อมูล เพื่อพัฒนาขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพให้มีประสิทธิภาพดีขึ้น รองรับข้อมูลจำนวนมากได้ โดยเน้นพัฒนาด้านความเร็วเปรียบเทียบกับโปรแกรม ClustalW

วัตถุประสงค์

1. ศึกษา และเปรียบเทียบขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพเพื่อเพิ่มประสิทธิภาพของการเทียบเรียง
2. พัฒนาขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพ โดยใช้เทคนิคการพิจารณาความถี่ส่วนย่อยของลำดับเข้าช่วย เพื่อให้การเทียบเรียงกลุ่มลำดับข้อมูลมีความเร็วเพิ่มขึ้น และสามารถรองรับการเทียบเรียงกลุ่มลำดับข้อมูลขนาดใหญ่ได้ โดยคำนึงถึงความถูกต้องของการเทียบเรียงด้วย

การตรวจเอกสาร

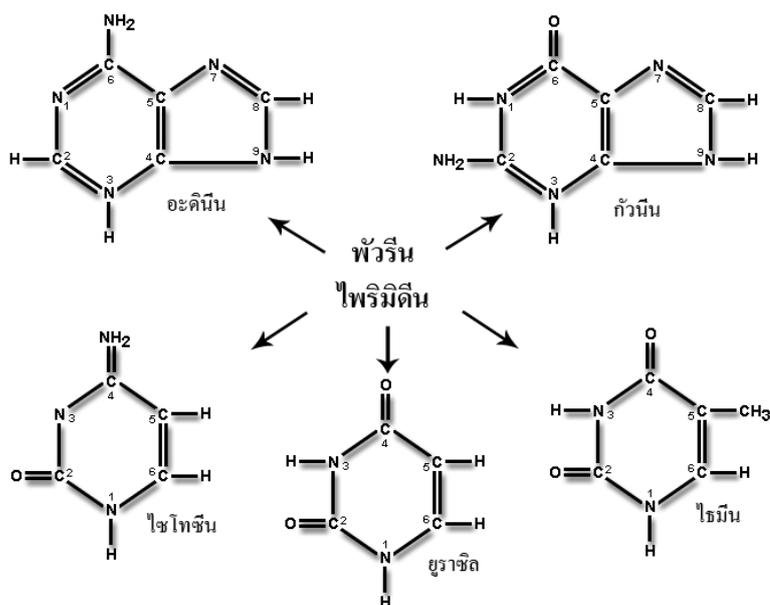
ความรู้เบื้องต้นเกี่ยวกับข้อมูลชีวภาพ

ในส่วนนี้กล่าวถึงความรู้พื้นฐานของข้อมูลชีวภาพที่ใช้ในการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพที่ใช้ในการเทียบเรียงแบ่งได้เป็น 2 กลุ่ม คือ ข้อมูลกรดนิวคลีอิก ได้แก่ ดีเอ็นเอ และ อาร์เอ็นเอ และข้อมูลโปรตีน

1. กรดนิวคลีอิก

กรดนิวคลีอิก(Nucleic acid) (ไพศาล, 2539; อมรา, 2540; พจน์ และคณะ, 2543) พบอยู่ในนิวเคลียส(Nucleus), ไซโทพลาซึม(Cytoplasm), และในออร์แกเนลล์(organelle) บางชนิดของเซลล์ สามารถแบ่งได้เป็น 2 ชนิดด้วยกันคือ ดีเอ็นเอ(Deoxyribonucleic acid: DNA) และอาร์เอ็นเอ(Ribonucleic acid: RNA) กรดนิวคลีอิกแต่ละชนิดประกอบด้วยเส้นสายที่เกิดจากการเชื่อมต่อของนิวคลีโอไทด์(nucleotide) หลายชนิด ในแต่ละนิวคลีโอไทด์ยังประกอบด้วยเบสที่มีไนโตรเจน(nitrogenous base) น้ำตาลที่มีคาร์บอน 5 อะตอม(pentose) และกรดฟอสฟอริก(phosphoric acid)

สารประกอบเบสที่มีไนโตรเจนที่พบในนิวคลีโอไทด์ เป็นสารที่มีวงแหวนคาร์บอน-ไนโตรเจนจำนวน 1 หรือ 2 วง เบสในนิวคลีโอไทด์แบ่งออกได้เป็น 2 กลุ่มตามจำนวนของวงแหวน คือ เบสที่มีวงแหวนเดี่ยว หรือไพริมิดีน(Pyrimidine) มีอยู่ 3 ชนิด ได้แก่ ไธมีน(Thymine: T), ไซโทซีน(Cytosine: C), และยูราซิล(Uracil: U) ส่วนเบสที่มี 2 วงแหวน หรือพิวรีน(Purine) มีอยู่ 2 ชนิด ได้แก่ อะดีนีน(Adenine: A) และกัวนีน(Guanine: G) เบสในนิวคลีโอไทด์มีโครงสร้างโมเลกุลดังภาพที่ 1

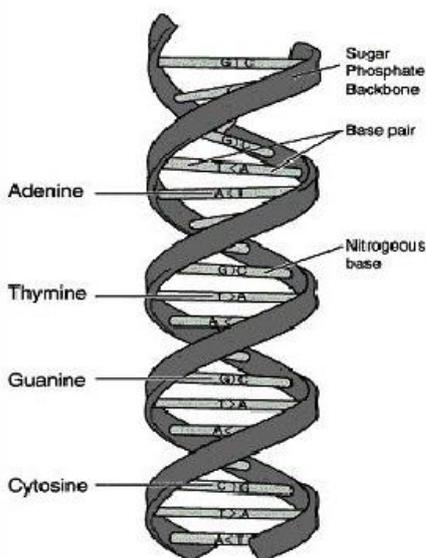


ภาพที่ 1 โครงสร้างโมเลกุลของเบสอันเป็นองค์ประกอบของนิวคลีโอไทด์

1.1 ดีเอ็นเอ เป็นสารชีวโมเลกุลที่มีขนาดใหญ่มาก เกิดจากการเรียงต่อกันของนิวคลีโอไทด์ บางโมเลกุลของดีเอ็นเออาจมีความยาวถึง 200,000 นิวคลีโอไทด์ ในเซลล์ของมนุษย์สามารถพบดีเอ็นเอได้ทั้งในนิวเคลียส และในไมโทคอนเดรีย สำหรับในพืชนอกจากจะพบดีเอ็นเอในนิวเคลียสแล้ว ยังพบดีเอ็นเอได้อีกในคลอโรพลาสต์อีกด้วย ดีเอ็นเอทำหน้าที่เป็นสารพันธุกรรม(Genetic Material) โดยจัดเรียงตามข้อกำหนดจำเพาะ เปรียบเสมือนเป็นแบบพิมพ์ชีวิต

โมเลกุลของดีเอ็นเอประกอบด้วยนิวคลีโอไทด์ 4 ชนิด ได้แก่ อะดีนีน, กัวนีน, ไทมิน, และไซโทซีน โมเลกุลของดีเอ็นเออยู่ในสภาพที่เป็นเกลียวคู่ขนาน(double helix) คล้ายกับบันไดเวียนดังภาพที่ 2 โดยมีน้ำตาลและกรดฟอสฟอริกเป็นราวบันไดทั้งสองข้าง พันธะไฮโดรเจนซึ่งยึดระหว่างพิวรีน และไพริมิดีนเป็นขั้นบันได โดยที่บันไดดังกล่าวหมุนเป็นเกลียวอย่างมีระเบียบและลำดับการเรียงตัวของเบสมีความเฉพาะเจาะจง ซึ่งเป็นหัวใจสำคัญในการควบคุมลักษณะทางพันธุกรรมและจะถูกถ่ายทอดต่อไป

ข้อมูลของดีเอ็นเอถูกเก็บในรูปแบบของตัวอักษร 4 ตัวคือ A G C และ T แทนนิวคลีโอไทด์ทั้ง 4 ชนิด และตัวอักษร N แทนเบสที่ไม่รู้จัก โดยที่ตัวอักษรจะเรียงเป็นสายอักษรตามลำดับของโครงสร้างโมเลกุลของดีเอ็นเอ



ภาพที่ 2 โครงสร้างของดีเอ็นเอ เป็นแบบเกลียวคู่ขนาน (double helix)

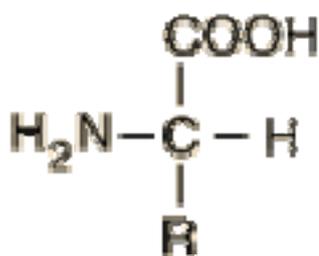
1.2 อาร์เอ็นเอ สามารถพบได้ในนิวเคลียส, ไซโทพลาซึม และไมโทคอนเดรีย มีหน้าที่เกี่ยวข้องกับการสังเคราะห์โปรตีน (protein synthesis) โครงสร้างโมเลกุลของอาร์เอ็นเอมีลักษณะคล้ายกับดีเอ็นเอ คือเป็นสายของนิวคลีโอไทด์ แต่แตกต่างจากดีเอ็นเอคือเป็นสายเดี่ยว และโมเลกุลของอาร์เอ็นเอจะประกอบด้วยนิวคลีโอไทด์ 4 ชนิด คือ อะดีนีน, กัวนีน, ไซโทซีน, และ ยูราซิล การเก็บข้อมูลของอาร์เอ็นเอมีรูปแบบเหมือนกับดีเอ็นเอ คือเป็นสายของตัวอักษร สายอักษรของข้อมูลอาร์เอ็นเอจะประกอบด้วยตัวอักษร 4 ตัว คือ A G C และ U

2. โปรตีน

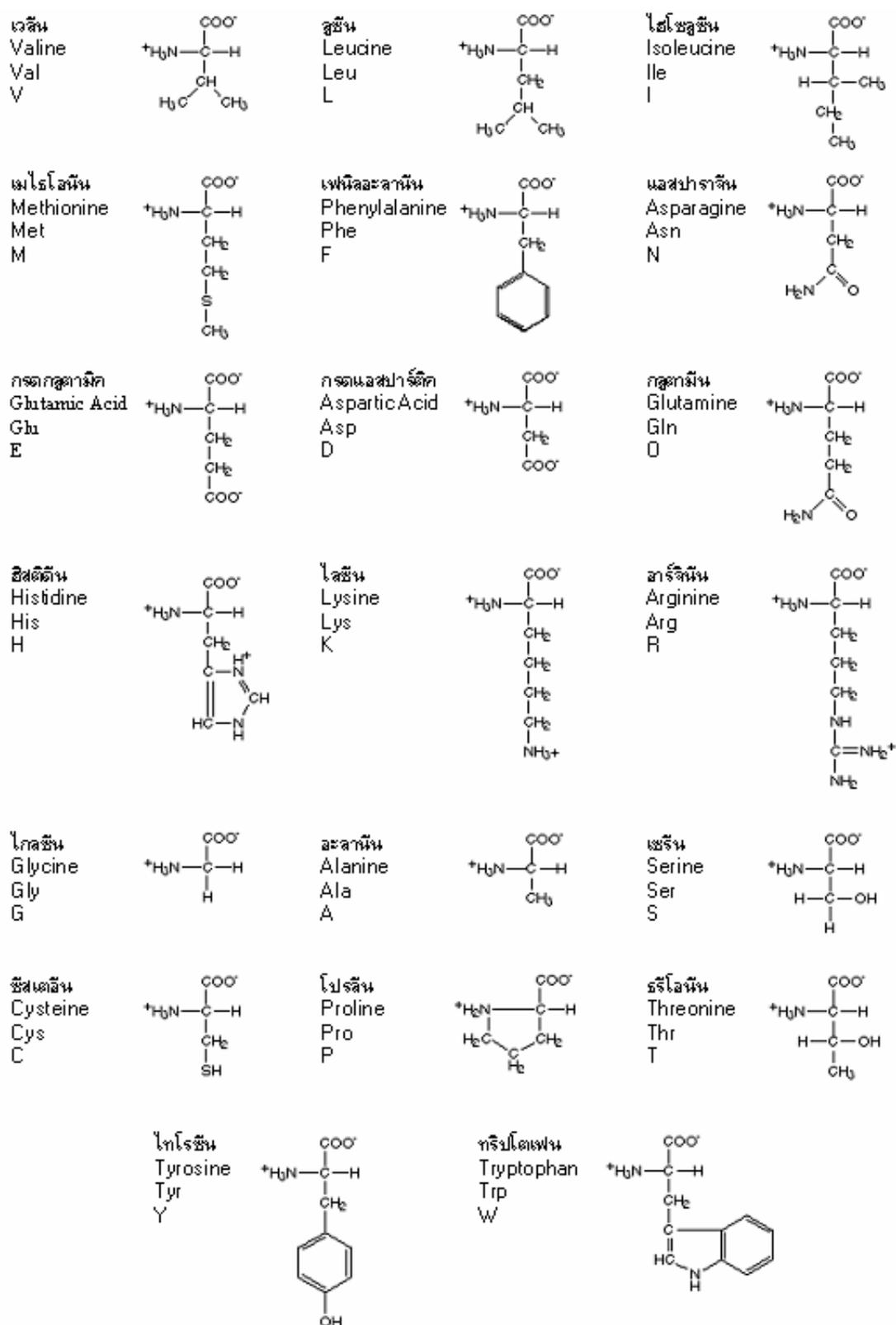
โปรตีน (ไพศาล, 2539; พจน์ และคณะ, 2543) จัดเป็นส่วนประกอบที่สำคัญของเซลล์ และเป็นส่วนประกอบของเอนไซม์ทุกชนิด โปรตีนและเอนไซม์จัดว่ามีความสำคัญในการก่อให้เกิดลักษณะต่าง ๆ ในสิ่งมีชีวิต โปรตีนจัดเป็นโมเลกุลขนาดใหญ่ โปรตีนแต่ละชนิดประกอบด้วยเส้นโพลีเปปไทด์ (Polypeptide) จำนวน 1 เส้น หรือมากกว่า 1 เส้น โพลีเปปไทด์ประกอบด้วยกรดอะมิโน (Amino acid) หลายชนิดมาเรียงต่อกัน กรดอะมิโนเป็นกรดอินทรีย์ที่ประกอบด้วยหมู่คาร์บอกซิลิก (Carboxylic group, $-\text{COOH}$), หมู่อะมิโน (Amino group, $-\text{NH}_2$), อะตอมไฮโดรเจน, และหมู่ R (side chain) ต่ออยู่กับอะตอมคาร์บอนในตำแหน่งแอลฟา (α -carbon) ดังภาพที่ 3 เนื่องจากหมู่ทั้งสี่ที่ต่อกับอะตอมคาร์บอนนี้ไม่เหมือนกัน คาร์บอนนี้จึงเป็นคาร์บอนชนิดอสมมาตร (Asymmetric carbon) ดังนั้นกรดอะมิโนทุกตัว (ยกเว้นกรดอะมิโนไกลซีนซึ่งมีหมู่ R เป็นอะตอมไฮโดรเจน) จะ

มีสเตอริโอไอโซเมอร์ได้ 2 ชนิด คือ D- และ L- กรดอะมิโนที่พบในร่างกายมนุษย์เป็นชนิด L- กรดอะมิโนมีอยู่มากกว่า 80 ชนิด แต่ชนิดที่พบทั่วไปในโปรตีนชนิดต่าง ๆ และถือว่ามีค่าสำคัญ มีอยู่เพียง 20 ชนิดเท่านั้น ซึ่งมีโครงสร้างโมเลกุล ชื่อ และชื่อย่อดังภาพที่ 4

ข้อมูลของโปรตีนถูกเก็บอยู่ในรูปแบบตัวอักษร 20 ตัวคือ ARND CQEGHILKMF PSTWY และ V แทนกรดอะมิโนทั้ง 20 ชนิด ตามชื่อย่อแบบหนึ่งตัวอักษร โดยตัวอักษรจะเรียงตัวเป็นสายอักษรตามการเรียงตัวของกรดอะมิโนภายในเส้นของโพลีเปปไทด์



ภาพที่ 3 โครงสร้างทั่วไปทางเคมีของกรดอะมิโน



ภาพที่ 4 โครงสร้างโมเลกุล และชื่อของกรดอะมิโนทั้ง 20 ชนิด

รูปแบบการจัดเก็บข้อมูลชีวภาพ

รูปแบบการจัดเก็บข้อมูลชีวภาพแสดงในงานวิจัยนี้จัดเก็บอยู่ในรูปแบบ FASTA (David, 2001; วีรุฒิ, 2545) ดังภาพที่ 5 โดยมี 3 ส่วน คือส่วนที่ 1 มีเครื่องหมาย > ตามด้วยชื่อและแหล่งที่มาของข้อมูลแต่ละสายอักษร ส่วนที่ 2 เป็นข้อมูลลำดับนิวคลีโอไทด์หรือกรออะมิโน ซึ่งใช้อักษรตัวเดียวเป็นสัญลักษณ์ ส่วนที่ 3 เป็นเครื่องหมาย * แสดงถึงจุดสิ้นสุดของข้อมูล โดยส่วนนี้จะมีหรือไม่มีก็ได้ เครื่องหมาย * อาจจำเป็นต้องใช้ในบางโปรแกรม

```
>YCZ2_YEAST protein in HMR 3' region
MKAVVIEDGKAVVKEGVPIPELEEGFV
GNPTDWAHIDYKVGPPQGSILGCDAAGQ
IVKLGPAVDPKDFSIGDYIYGFIHGSS
VRFPSNGAFAEYSAI STVVAYKSPNEL
KFLGEDVLPAGPVRSLGAATIPVSLT*
```

ภาพที่ 5 ตัวอย่างการจัดเก็บข้อมูลชีวภาพแบบ FASTA

ลักษณะการเปลี่ยนแปลงของข้อมูลชีวภาพ

การเปลี่ยนแปลงในดีเอ็นเอ หรือโปรตีน จะมีลักษณะเหมือนกันคือจะเกิดจากการแทรก การลบ หรือการเปลี่ยนของตัวอักษร ดังตัวอย่างภาพที่ 6 แสดงการเปลี่ยนของลำดับเบสในดีเอ็นเอ การเปลี่ยนแปลงของข้อมูลชีวภาพในลักษณะต่าง ๆ เหล่านี้ส่งผลให้ลักษณะของสิ่งมีชีวิต แตกต่างไปจากเดิมด้วย(ไพศาล, 2539; อมรา, 2540; วีรุฒิ, 2545)

- 1) AAGGCAAACCTACTGGTCTTATG (สายปกติ)
- 2) AAGGCAAATCTACTGGTCTTATG (เปลี่ยนคนละคู่เบสกัน)
- 3) AAGGCAAACCTACTGCTCTTATG (เปลี่ยนคู่เบสเดียวกัน)
- 4) AAGGCAA^{ACCTA}ACTGCTCTTATG (การขาดไป)
- 5) AAGGCAAACCTACTAAAGCGCTCTTATG (การเพิ่มเข้ามา)
- 6) AAGGTTTGCCTACTGGTCTTATG (การอินเวอร์ชัน)

ภาพที่ 6 รูปแบบการเปลี่ยนแปลงของเบสในดีเอ็นเอ โดย 1) คือลำดับเบสของดีเอ็นเอสายปกติ 2) คือการเปลี่ยนจาก C ไป T เป็นการเปลี่ยนของคนละคู่เบสกัน(Transition) 3) คือการเปลี่ยนจาก G ไป C เป็นการเปลี่ยนกับคู่คอมพลีเมนต์กัน (Transversion) 4) คือการขาดไปของช่วงลำดับเบส ACCTA 5) คือการเพิ่มเข้ามาของช่วงลำดับเบส AAAGC 6) คือการอินเวอร์ชัน จาก 5-GCAAAC-3 เปลี่ยนเป็น 5-GTTTGC-3

การเทียบเรียงข้อมูล

การเทียบเรียงข้อมูล(Sequence Alignment) (Crochemore, 2003) เป็นวิธีสำหรับการเปรียบเทียบข้อมูล เพื่อระบุส่วนที่คล้าย และส่วนที่แตกต่าง ของข้อมูลที่ต้องการเปรียบเทียบ การเทียบเรียงจะพิจารณาจากความคล้ายกันของอักขระในแต่ละตำแหน่งของลำดับข้อมูล ซึ่งค่าความคล้ายกันจะคำนวณจากรูปแบบการให้คะแนนที่เหมาะสม

รูปแบบเบื้องต้นในการเทียบเรียงข้อมูลจำนวน 2 ลำดับเริ่มจากเขียนลำดับแรกไว้บรรทัดบนของที่ 2 แล้วพิจารณาส่วนย่อยที่มีความเหมือนกันระหว่าง 2 ลำดับ ในแต่ละลำดับสามารถแยกส่วนย่อยออกจากกันได้ด้วยช่องว่าง ซึ่งโดยปกติแล้วการแทรกช่องว่างจะไม่แทรกในตำแหน่งที่ตรงกันของทั้งสองลำดับ ผลลัพธ์สุดท้ายที่ได้หลังจากเทียบเรียงแล้วลำดับที่เทียบเรียงจะมีความยาวเท่ากัน ตัวอย่างการเทียบเรียงลำดับ A = “ACAAGACAGCGT” และลำดับ B = “AGAACAAGGCGT” เป็นดังภาพที่ 7

```

A = ACAAGACAG-CGT
    | | | | |
B = AGAACA-AGGCGT

```

ภาพที่ 7 ตัวอย่างการเทียบเรียงลำดับ 2 ลำดับ

จากตัวอย่างภาพที่ 7 จะเห็นได้ว่าผลการเทียบเรียงมีอักขระที่เหมือนกันในตำแหน่งเดียวกันทั้งหมด 9 ตัว อย่างไรก็ตามการเทียบเรียงให้อักขระที่ต่างกันอยู่ในตำแหน่งเดียวกันสามารถเกิดขึ้นได้ ดังเช่นในตัวอย่างตำแหน่งที่ 2 อักขระ C ของลำดับ A ถูกจัดวางในตำแหน่งที่ตรงกับอักขระ G ของลำดับ B และตำแหน่งที่ 4 อักขระ G ของลำดับ A ก็ถูกจัดวางในตำแหน่งที่ตรงกับอักขระ C ของลำดับ B เช่นเดียวกัน นอกจากการเทียบเรียงอักขระที่เหมือนและไม่เหมือนแล้ว การแทรกช่องว่างภายในลำดับก็สามารถทำได้เช่นเดียวกัน เช่นในตำแหน่งที่ 10 ของลำดับ A การแทรกช่องว่างทำให้ได้การเทียบเรียงที่ดีสำหรับอักขระ 3 ตัวสุดท้าย

หากมองในมุมมองของการเปลี่ยนแปลงของข้อมูลแล้ว ผลของการเทียบเรียงข้อมูลสามารถแสดงให้เห็นถึงรูปแบบการเปลี่ยนแปลงจากลำดับหนึ่งไปเป็นอีกลำดับหนึ่งได้ โดยในตำแหน่งที่อักขระไม่ตรงกันถือว่าการแทนที่(substitution) ของอักขระ ในตำแหน่งที่เป็นช่องว่างในลำดับแรกถือว่าการแทรก(insertion) อักขระที่ลำดับที่ 2 และในตำแหน่งที่เป็นช่องว่างในลำดับที่ 2 ถือว่าการลบ(deletion) อักขระที่ลำดับแรกออก จากตัวอย่างการเทียบเรียงในภาพที่ 7 สามารถสรุปได้ว่าการเปลี่ยนแปลงจากลำดับ A ไปเป็นลำดับ B มี 4 ขั้นตอน คือ 1) แทนที่อักขระ C ด้วย G ในตำแหน่งที่ 2; 3) แทนที่อักขระ G ด้วย C ในตำแหน่งที่ 5; 4) ลบอักขระ C ในตำแหน่งที่ 7; และ 4) แทรกอักขระ G ในตำแหน่งที่ 10

การเทียบเรียงข้อมูลนั้นสามารถมีผลการเทียบเรียงได้มากมายหลายรูปแบบ แต่ผลการเทียบเรียงที่สนใจนั้น คือผลที่เทียบเรียงข้อมูลให้มีความคล้ายกันมากที่สุด วิธีการวัดผลความคล้ายกันของการเทียบเรียงข้อมูลสามารถคำนวณได้จากค่าคะแนนความคล้ายของแต่ละคู่อักขระซึ่งมีหลายรูปแบบ โดยปกติแล้วค่าคะแนนความคล้ายจะให้คะแนนที่ดีสำหรับการเทียบเรียงอักขระเดียวกัน และให้คะแนนที่ไม่ดี หรือหักคะแนนสำหรับการเทียบเรียงอักขระที่ต่างกันหรือช่องว่าง คะแนนรวมของการเทียบเรียงจะคำนวณจากผลรวมคะแนนในแต่ละตำแหน่งของผลการเทียบเรียง ตัวอย่างเช่นหากวิธีการให้คะแนนอยู่ในรูปแบบให้คะแนน +1 สำหรับการเทียบเรียงอักขระเดียวกัน

และให้คะแนน -1 สำหรับการเทียบเรียงอักขระต่างกันหรือช่องว่าง จากตัวอย่างการเทียบเรียงในภาพที่ 7 มีการเทียบเรียงอักขระเดียวกัน 9 ตัว และมีการเทียบเรียงอักขระที่ต่างกัน หรือช่องว่าง 4 ตัว ดังนั้นคะแนนความคล้ายกันจะมีค่าเท่ากับ 5 เป็นต้น

เมื่อพิจารณาถึงจำนวนข้อมูลที่ต้องการเทียบเรียง การเทียบเรียงข้อมูลสามารถแบ่งได้เป็น 2 แบบ คือ การเทียบเรียงคู่ลำดับข้อมูล(Pairwise Alignment) และการเทียบเรียงกลุ่มลำดับข้อมูล(Multiple Alignment)

การเทียบเรียงคู่ลำดับข้อมูล คือ การเทียบเรียงลำดับข้อมูลจำนวน 2 ลำดับ เพื่อหาความคล้าย, ความแตกต่าง หรือการเปลี่ยนแปลงระหว่างข้อมูล และนอกจากนั้นยังสามารถประยุกต์ไปใช้กับการค้นหากลุ่มลำดับที่ใกล้เคียง(Homologous sequence) จากฐานข้อมูล สำหรับลำดับที่ต้องการได้อีกด้วย

การเทียบเรียงกลุ่มลำดับข้อมูล(Gusfield, 1999) คือการเทียบเรียงลำดับข้อมูลจำนวนตั้งแต่ 3 ลำดับขึ้นไป เพื่อหาความสัมพันธ์โดยรวมของข้อมูลทั้งหมดที่เทียบเรียง โดยจุดประสงค์หลักคือการหาส่วนที่คล้ายกันจากทุกลำดับที่เทียบเรียง ซึ่งเป็นวิธีการที่ขยายมาจากการเทียบเรียงคู่ลำดับข้อมูล

เมื่อพิจารณาถึงช่วงของลำดับที่ต้องการเทียบเรียงข้อมูล การเทียบเรียงข้อมูลสามารถแบ่งได้เป็น 2 แบบ คือ การเทียบเรียงแบบโกลบอล(Global Alignment) และการเทียบเรียงแบบโลคอล(Local Alignment)

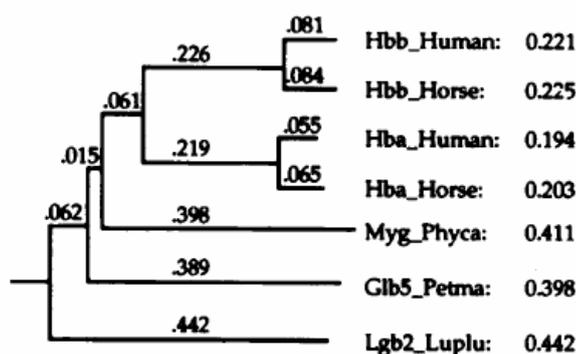
การเทียบเรียงแบบโกลบอล(Gusfield, 1999:216-217) เป็นการหาการเทียบเรียงที่ดีที่สุดของทั้งสายลำดับข้อมูล โดยการเทียบเรียงแบบโกลบอลมีนิยาม คือ การเทียบเรียงแบบโกลบอลของลำดับข้อมูล S1 และ S2 จำนวน 2 ลำดับ ทำได้โดยแทรกช่องว่างภายในลำดับ หรือที่จุดสิ้นสุดของลำดับ ของทั้งลำดับ S1 และ S2 หลังจากนั้นจัดวางผลการเทียบเรียงของทั้ง 2 ลำดับ โดยให้ลำดับหนึ่งอยู่เหนืออีกลำดับหนึ่ง ผลที่ได้จากการเทียบเรียงคือทุกๆอักขระหรือช่องว่างในแต่ละลำดับต้องตรงข้ามกับอักขระ หรือช่องว่างของอีกลำดับเพียงตัวเดียวเท่านั้น

จากนิยามของการเทียบเรียงแบบโกลบอล คือผลที่ได้จากการเทียบเรียง ทุกลำดับจะมีความยาวเท่ากัน และแสดงให้เห็นถึงการจัดเรียงอักษรของทั้งลำดับข้อมูลให้มีความคล้ายกันมากที่สุด ซึ่งผลที่ได้จะแตกต่างกันไปขึ้นอยู่กับวิธีการ และค่าการให้คะแนนที่เลือกใช้

การเทียบเรียงแบบโลคอล(Gusfield, 1999:230-232) เป็นการหาส่วนย่อยของลำดับที่ให้ผลการเทียบเรียงที่ดีที่สุด โดยการเทียบเรียงแบบโลคอลมีนิยาม คือ การเทียบเรียงแบบโลคอลของลำดับข้อมูล S1 และ S2 คือการหาส่วนย่อยของลำดับข้อมูล S1 และ S2 ที่ให้ผลการเทียบเรียงแบบโกลบอลระหว่างส่วนย่อยสูงที่สุด และดีกว่าทุกคู่ของส่วนย่อยอื่นที่ได้จาก S1 และ S2

ประโยชน์ของการเทียบเรียงกลุ่มข้อมูลชีวภาพ

การเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพเป็นวิธีการที่ทำให้ทราบถึงความสัมพันธ์ภายในโครงสร้างของสายโปรตีน ซึ่งสามารถนำไปประยุกต์ใช้กับงานทางด้านชีววิทยาได้หลากหลาย งานทางด้านชีววิทยาที่นิยมใช้วิธีการเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพมีหลักๆ 3 แนว คือ วิเคราะห์วิวัฒนาการของสิ่งมีชีวิต(Phylogenetic analyses) ดังตัวอย่างภาพที่ 8 ระบุส่วนจำเพาะ (Conserved motifs and domains) ของกลุ่มข้อมูลดังตัวอย่างภาพที่ 9 และทำนายโครงสร้างทุติยภูมิ และตติยภูมิ(Secondary and tertiary structure prediction) ของข้อมูลโปรตีน(Notredame, 2002)



ภาพที่ 8 ตัวอย่างแผนภูมิด้านไม่แสดงความสัมพันธ์ทางวิวัฒนาการที่วิเคราะห์ได้จากการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพ

			*	*	*	*	*
.....	.MATQDSEVA	LVTGATSGIG	LEIARRLGKE	GLRVFV.	CAR		
MTTATATATA	TPGTAAKPVA	LVTGATSGIG	LAIARRLAAL	GARTFL.	CAR		
.....	MTTSATPRTA	LVTGGTSGIG	LAVVKTLAAR	GLRVFL.	CAR		
.....	.MNFEGKIA	LVTGASRGIG	RAIAETLAAR	GGKVIQ.	...		
.....	.MFELTGRKA	LVTGASGAIG	GAIARVLHAQ	GAIVGL.	...		
.....MTQRIA	YVTGGMGGIG	TAICQRLAKD	GFRVVA.	GCG		
.....	MYTDLKDKVV	VITGGSTGLG	RAMAVRFGQE	EAKVVI.	NYI		
			*	*	*	*	*

ภาพที่ 9 ตัวอย่างผลของการเทียบเรียงกลุ่มลำดับข้อมูลเพื่อระบุส่วนจำเพาะของกลุ่มลำดับข้อมูล โดยส่วนที่มีเครื่องหมาย * คือส่วนจำเพาะของกลุ่มลำดับข้อมูล

ตารางค่าคะแนน

ตารางค่าคะแนนในการเทียบเรียงกลุ่มลำดับข้อมูล (วิรุฒิ, 2545) เป็นค่าคะแนนที่บอกความเหมือนของแต่ละตัวอักษรของข้อมูลชีวภาพ สร้างค่าคะแนนโดยการหาทางเคมีโดยแต่ละตารางมีค่าคะแนนที่ต่างกัน การเลือกใช้ตารางค่าคะแนนก็มีผลต่อความถูกต้องของการเทียบเรียง โดยงานวิจัยนี้ใช้ตารางค่าคะแนน Gonnet250 ซึ่งเป็นตารางค่าคะแนนที่ปกติโปรแกรม ClustalW (Thomson et al., 1994) ใช้ นอกจากนี้ยังสามารถเลือกใช้ตารางค่าคะแนนได้หลายแบบซึ่งเปลี่ยนได้โดยผู้ใช้ ได้แก่ BLOSUM PAM และ Gonnet ซึ่งเป็นตารางค่าคะแนนที่นิยมใช้

1. BLOSUM

BLOSUM เป็นตารางค่าคะแนนที่คิดค้นโดย Henikoff ในปี ค.ศ. 1992 ซึ่งเหมาะกับการค้นหาความเหมือน (similarity searches) เหมาะกับโปรแกรมในการค้นหา ได้แก่ FASTA และ BLAST โดยที่ตารางค่าคะแนน BLOSUM มีหลายระดับ แต่ละระดับมีค่ากำหนดที่ต่างกัน เช่น BLOSUM62 กำหนดจากชุดโปรตีนที่มีความเหมือนกันมากกว่า 62% เป็นต้น

2. PAM

PAM เป็นตารางค่าคะแนนโดย Dayhoff ซึ่งเป็นที่นิยมตั้งแต่ช่วงปี ค.ศ. 1978 โดยมีพื้นฐานจากร้อยละของการยอมรับการเปลี่ยนรูปของสายอักษร (mutation) และการอ้างถึงการเปลี่ยนในโปรตีน โดยการทดสอบความถูกต้อง Dayhoff พบว่าตารางค่าคะแนน PAM เหมาะสำหรับการหาค่าความห่างของความสัมพันธ์

3. Gonnet

Gonnet เป็นตารางค่าคะแนนโดย Cohen และ Benner ในปี ค.ศ. 1992 ซึ่งตารางค่าคะแนนนี้พัฒนามาจากตารางค่าคะแนน PAM โดยตารางค่าคะแนน Gonnet จะอิงชุดข้อมูลที่มีขนาดใหญ่

เทคนิคที่ใช้ในการเทียบเรียงกลุ่มข้อมูล

ในส่วนนี้จะกล่าวถึงเทคนิคที่ใช้ในการเทียบเรียงกลุ่มลำดับข้อมูล ซึ่งจะกล่าวถึง 2 เทคนิคหลัก คือ การโปรแกรมแบบพลวัต และฮิวริสติก

1. การโปรแกรมแบบพลวัต(Dynamic Programming)

การโปรแกรมแบบพลวัต (Dynamic Programming) (Gusfield, 1999; Notredame, 2002; วรวิทย์, 2545) เป็นวิธีการที่ให้ค่าการเทียบเรียงที่ดีที่สุด โดยการคำนวณอ้างอิงจากรางค่าคะแนนที่เลือกใช้ ซึ่งวิธีการหาเริ่มจากการเลือกตารางค่าคะแนน และค่าของช่องว่างที่ต้องการใช้ และสร้างตารางการคำนวณ โดยให้ลำดับข้อมูลแรกกำกับแนวสดมภ์ และลำดับข้อมูลที่สองกำกับแนวแถวของตารางการคำนวณ และมีเงื่อนไขพื้นฐานตามสมการที่ 1 และ สมการที่ 2

$$F(i, 0) = d \quad (1)$$

$$F(0, j) = d \quad (2)$$

โดยที่ i และ j คือตำแหน่งของสดมภ์ และแถวในตารางการคำนวณ $F(i, j)$ คือค่าในตารางที่ได้จากการคำนวณแบบพลวัต และ d คือค่าของช่องว่างที่เลือกใช้

จากเงื่อนไขพื้นฐานตามสมการที่ 1 และ สมการที่ 2 การคำนวณต่อมาเป็นการเติมค่าในตารางการคำนวณให้ครบ โดยการเทียบเรียงแบบโกลบอลใช้สมการที่ 3 ส่วนการเทียบเรียงแบบโลคอลจะเริ่มต้นค่าพื้นฐานเป็นศูนย์ และใช้สมการที่ 4 หาค่าเพื่อใส่ในตารางการคำนวณวิธีการคำนวณค่าในตารางเป็นดังภาพที่ 10

กำหนดให้

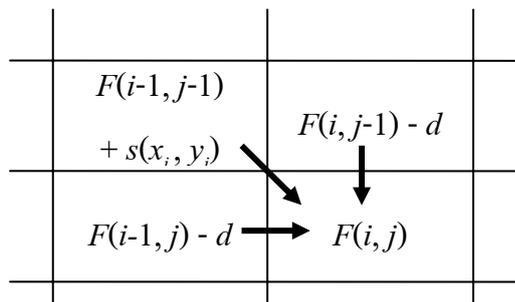
$F(i, j)$ เป็นค่าที่คำนวณได้ในตารางการคำนวณ

$s(x_i, y_j)$ เป็นค่าเปรียบเทียบตัวอักษรในตำแหน่งนั้น อ้างอิงจากตารางค่าคะแนน

d เป็นค่าช่องว่างที่เลือกใช้

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j), \\ F(i-1, j) - d, \\ F(i, j-1) - d. \end{cases} \quad (3)$$

$$F(i, j) = \max \begin{cases} 0, \\ F(i-1, j-1) + s(x_i, y_j), \\ F(i-1, j) - d, \\ F(i, j-1) - d. \end{cases} \quad (4)$$



ภาพที่ 10 การคำนวณค่าในตารางการคำนวณ

เมื่อคำนวณจนได้ค่าในตารางการคำนวณครบแล้ว ผลลัพธ์ของการเทียบเรียงข้อมูลได้จากการหาเส้นทางย้อนกลับจากมุมล่างขวาของตารางการคำนวณย้อนทางค่าที่คำนวณได้จนถึงตำแหน่งบนซ้ายของตารางก็จะได้การเทียบเรียงข้อมูลที่ให้ค่าคะแนนที่ดีที่สุด โดยหากเส้นทางการย้อนกลับเป็น $i-1$ และ $j-1$ หมายถึงอักขระตำแหน่ง i และ j ถูกเทียบเรียงอยู่ในตำแหน่งที่ตรงกัน หากการย้อนกลับเป็น $i-1$ และ j หมายถึงอักขระในสดมภ์ตำแหน่ง i ถูกเทียบเรียงคู่กับช่องว่าง และหากการย้อนกลับเป็น i และ $j-1$ หมายถึงอักขระในแถวตำแหน่งที่ j ถูกเทียบเรียงคู่กับช่องว่าง ซึ่งตัวอย่างการคำนวณแบบพลวัตเป็นดังตัวอย่างที่ 1

ตัวอย่างที่ 1 การเทียบเรียงแบบโกลบอล และการเทียบเรียงแบบโลคอล ของลำดับข้อมูล HEAGAWGHEE และ PAWHEAE โดยวิธีการโปรแกรมแบบพลวัต และใช้ตารางค่าคะแนน BLOSUM50 ดังตารางที่ 1 และใช้ค่าของช่องว่างคือ -8

ตารางที่ 1 ตารางค่าคะแนน BLOSUM50

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	5	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-3	-1	1	0	-3	-2	0
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1	-3
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2	-3
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3	-4
C	-1	-4	-2	-4	13	-3	-3	-3	-3	-2	-2	-3	-2	-2	-4	-1	-1	-5	-3	-1
Q	-1	1	0	0	-3	7	2	-2	1	-3	-2	2	0	-4	-1	0	-1	-1	-1	-3
E	-1	0	0	2	-3	2	6	-3	0	-4	-3	1	-2	-3	-1	-1	-1	-3	-2	-3
G	0	-3	0	-1	-3	-2	-3	8	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4
H	-2	0	1	-1	-3	1	0	-2	10	-4	-3	0	-1	-1	-2	-1	-2	-3	2	-4
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5	2	-3	2	0	-3	-3	-1	-3	-1	4
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5	-3	3	1	-4	-3	-1	-2	-1	1
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6	-2	-4	-1	0	-1	-3	-2	-3
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7	0	-3	-2	-1	-1	0	1
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8	-4	-3	-2	1	4	-1
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10	-1	-1	-4	-3	-3
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5	2	-4	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5	-3	-2	0
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-4	15	2	-3
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	-1
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5

การเทียบเรียงแบบโกลบอลเริ่มจากการสร้างตารางการคำนวณจากเงื่อนไขพื้นฐานตาม สมการ 1 และสมการที่ 2 และคำนวณค่าในแต่ละตำแหน่งของตารางการคำนวณโดยใช้สมการ 3 เมื่อได้คำนวณค่าในตารางจนครบจึงหาเส้นทางย้อนกลับจากมุมล่างขวาของตาราง ก็จะได้ผลการเทียบเรียงดังภาพที่ 11

	H	E	A	G	A	W	G	H	E	E	
P	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
A	-8	-2	-9	-17	-25	-33	-41	-49	-57	-65	-73
W	-16	-10	-3	-4	-12	-20	-28	-36	-44	-52	-60
H	-24	-18	-11	-6	-7	-15	-5	-13	-21	-29	-37
E	-32	-14	-18	-13	-8	-9	-13	-7	-3	-11	-19
A	-40	-22	-8	-16	-16	-9	-12	-15	-7	3	-5
E	-48	-30	-16	-3	-11	-11	-12	-12	-15	-5	2
E	-56	-38	-24	-11	-6	-12	-14	-15	-12	-9	1

HEAGAWGHE-E
-PA--W-HEAE

ภาพที่ 11 การเทียบเรียงแบบโกลบอลโดยใช้การโปรแกรมพลวัต และผลลัพธ์ที่ได้

การเทียบเรียงแบบโกลบอลเริ่มจากการสร้างตารางการคำนวณจากเงื่อนไขพื้นฐานเท่ากับศูนย์ และคำนวณค่าในแต่ละตำแหน่งของตารางการคำนวณโดยใช้สมการที่ 4 เมื่อคำนวณค่าในตารางจนครบแล้ว ผลการเทียบเรียงแบบโกลบอลหาได้จากการหาเส้นทางย้อนกลับจากจุดที่มีค่าคะแนนสูงที่สุดในตารางย้อนกลับไปยังกระทั่งมีค่าเป็นศูนย์ ดังภาพที่ 12 ค่าคะแนนสูงที่สุดในตารางคือ 28 และหาเส้นทางย้อนกลับไปยังจนถึงค่าศูนย์ก็จะได้ผลการเทียบเรียงแบบโกลบอลของทั้ง 2 ลำดับ

	H	E	A	G	A	W	G	H	E	E
P	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	0	0
W	0	0	0	5	0	5	0	0	0	0
H	0	0	0	0	2	0	20	12	4	0
E	0	10	2	0	0	0	12	18	22	14
A	0	2	16	8	0	0	4	10	18	28
E	0	0	8	21	13	5	0	4	10	20
E	0	0	6	13	18	12	4	0	4	16

AWGHE
AW-HE

ภาพที่ 12 การเทียบเรียงแบบโกลบอลโดยใช้การโปรแกรมพลวัต และผลลัพธ์ที่ได้

วิธีการเทียบเรียงลำดับข้อมูลโดยวิธีการโปรแกรมพลวัตเป็นวิธีที่รับประกันผลลัพธ์ที่ดีที่สุด โดยมีค่าความซับซ้อนของเวลา (Time Complexity) คือ $O(L^N)$ เมื่อ L คือความยาวเฉลี่ยของแต่ละลำดับ และ N คือจำนวนลำดับที่เข้าเทียบเรียง ดังนั้นหากใช้วิธีการโปรแกรมพลวัตเทียบเรียงลำดับข้อมูลที่มีความยาวของลำดับมาก หรือเทียบเรียงลำดับจำนวนมากจะใช้เวลาในการประมวลผลมาก วิธีการโปรแกรมพลวัตจึงไม่เหมาะสมกับการเทียบเรียงกลุ่มลำดับข้อมูลในระดับปริมาณมากกว่า 5 ลำดับที่มีความยาวของลำดับ 100 อักขระขึ้นไป (Gusfield, 1999)

2. ฮิวริสติก(Heuristic)

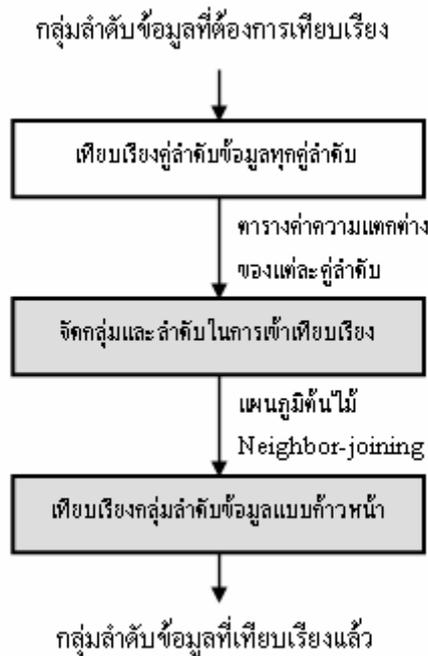
การเทียบเรียงลำดับข้อมูลโดยวิธีการฮิวริสติก เป็นการเทียบเรียงลำดับข้อมูลแบบโดยอ้างอิงหลักการประมาณ พยายามให้ผลการเทียบเรียงที่ได้มีคะแนนการเทียบเรียงใกล้เคียงกับค่าคะแนนที่ดีที่สุด วิธีในกลุ่มฮิวริสติกที่นิยมใช้ในการเทียบเรียงลำดับข้อมูลคือ การเทียบเรียงแบบก้าวหน้า(Progressive Alignment) (Notredame, 2002) ซึ่งเป็นวิธีที่เรียบง่าย และมีประสิทธิภาพสำหรับการเทียบเรียงกลุ่มลำดับข้อมูลโดยใช้เวลา และหน่วยความจำในการประมวลผลไม่มาก ทำให้วิธีนี้เป็นที่นิยมในการใช้งานในการเทียบเรียงปัจจุบัน เนื่องจากใช้เวลาในการประมวลผลน้อย และรองรับการเทียบเรียงข้อมูลขนาดใหญ่ หรือปริมาณมากได้(Thomson et al., 1999) โปรแกรมที่ใช้วิธีการในกลุ่มของฮิวริสติกและเป็นที่ยอมรับได้แก่ Pileup, MultiAlign, ClustalW และ ClustalX ซึ่งโปรแกรม ClustalX เป็น โปรแกรม ClustalW ที่แสดงผลเป็นแบบกราฟฟิก

โครงสร้างการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW

โปรแกรม ClustalW เป็นโปรแกรมเทียบเรียงกลุ่มลำดับข้อมูลโดยใช้วิธีการเทียบเรียงข้อมูลแบบก้าวหน้าที่ได้รับความนิยมเพราะเป็นโปรแกรมไม่คิดมูลค่า มีความถูกต้องสูง ใช้กับลำดับที่มากได้ และใช้ได้กับหลายระบบปฏิบัติการ ได้แก่ SUN Solaris, IRIX5.3 บนเครื่อง Silicon Graphics, Digital UNIX บน DEC Stations, Microsoft Windows(32 bit) สำหรับเครื่องคอมพิวเตอร์ส่วนบุคคล, Linux สำหรับเครื่องคอมพิวเตอร์ส่วนบุคคล และ Macintosh PowerMac (วีรวุฒิ, 2545)

โปรแกรม ClustalW สามารถเทียบเรียงได้ทั้งดีเอ็นเอ และโปรตีน โดยทุกลำดับที่ต้องการเทียบเรียงต้องอยู่ในไฟล์เดียวกัน โปรแกรม ClustalW รับไฟล์ได้ 7 รูปแบบ คือ NBRF/PIR, EMBL/SWISSPROT, Pearson (Fasta), Clustal (*.aln), GCG/MSF(Pileup), GCG9/RSF และ

GDE (Thomson *et al.*, 1994) การเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW แบ่งเป็น 3 ขั้นตอนดังภาพที่ 13



ภาพที่ 13 ขั้นตอนในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW

ขั้นตอนแรกเป็นการสร้างตารางค่าความแตกต่างของแต่ละคู่ลำดับข้อมูล (Distance matrix) เพื่อใช้เป็นค่าคะแนนสำหรับจัดลำดับการเทียบเรียงกลุ่มลำดับข้อมูลแบบก้าวหน้า การหาค่าคะแนนความแตกต่างนั้นใช้วิธีเทียบเรียงคู่ลำดับข้อมูลแต่ละคู่ โดยใช้เทคนิคการโปรแกรมพลวัตในการเทียบเรียงแต่ละคู่ลำดับข้อมูล เพื่อคำนวณหาค่าความคล้ายกันของแต่ละคู่ข้อมูล และการคำนวณค่าคะแนนความแตกต่างของแต่ละคู่ข้อมูล ($D_{ClustalW}(x, y)$) จะเป็นดังสมการที่ 5

$$D_{ClustalW}(x, y) = 1 - S_{ClustalW}(x, y) \quad (5)$$

โดยที่ $D_{ClustalW}(x, y)$ คือค่าคะแนนความแตกต่างของคู่ลำดับข้อมูล x และ y และ $S_{ClustalW}(x, y)$ คือค่าคะแนนความคล้ายกันของคู่ลำดับข้อมูล x และ y

การทำงานในขั้นตอนแรกใช้วิธีการโปรแกรมพลวัตที่ละคู่ลำดับ ซึ่งมีความซับซ้อนของเวลา $O(L^2)$ เมื่อ L คือความยาวของลำดับ และต้องเทียบเรียงคู่ลำดับให้ครบทุกคู่ลำดับ ซึ่งมีจำนวนการคำนวณทั้งหมด $N(N-1)/2$ ครั้ง หรือประมาณ $O(N^2)$ เมื่อ N คือจำนวนของลำดับข้อมูลที่ต้องการเทียบเรียง ดังนั้นในขั้นตอนนี้มีความซับซ้อนของเวลาทั้งหมด $O(N^2L^2)$

ในขั้นตอนที่ 2 เป็นการนำตารางค่าความแตกต่างของแต่ละคู่ลำดับข้อมูลมาสร้างเป็นแผนภูมิต้นไม้โดยใช้วิธีเนเบอร์-จอยนิง (neighbor-joining) (Saitou, 1987) เพื่อจัดกลุ่ม และจัดลำดับการเข้าเทียบเรียงกลุ่มลำดับข้อมูลแบบก้านหน้าสำหรับแต่ละลำดับข้อมูล

ขั้นตอนสุดท้ายเป็นการเทียบเรียงกลุ่มลำดับข้อมูลแบบก้านหน้าโดยลำดับการเข้าเทียบเรียงอ้างอิงจากแผนภูมิต้นไม้ที่สร้างในขั้นตอนที่ 2 การเทียบเรียงนั้นเริ่มจากคู่ลำดับข้อมูลที่มีค่าคะแนนความแตกต่างกันน้อยที่สุด หลังจากนั้นจึงนำลำดับข้อมูลที่มีค่าคะแนนความแตกต่างมากกว่าเข้ามารวมเทียบเรียงกับกลุ่มลำดับข้อมูลที่เทียบเรียงไปก่อนหน้านี้ ทำจนครบทุกลำดับก็จะได้ผลลัพธ์การเทียบเรียงกลุ่มลำดับข้อมูล

ปัญหาของการเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพ

การเทียบเรียงกลุ่มลำดับข้อมูลทั้งวิธีการโปรแกรมแบบพลวัต และวิธีในกลุ่มฮิวริสติกนั้นได้ผลที่มีความถูกต้องในระดับที่ยอมรับได้ แต่เมื่อนำมาเทียบเรียงกลุ่มลำดับที่มีความยาวของข้อมูลมาก หรือมีจำนวนมากก็จะเกิดปัญหาทางด้านเวลาในการประมวลผลได้

วิธีการโปรแกรมพลวัตเมื่อนำมาใช้กับการเทียบเรียงคู่ลำดับจะมีความซับซ้อนของเวลา คือ $O(L^2)$ และเมื่อใช้กับการเทียบเรียงกลุ่มลำดับข้อมูลจะมีค่าความซับซ้อนของเวลา คือ $O(L^N)$ เมื่อ L คือความยาวเฉลี่ยของแต่ละลำดับ และ N คือจำนวนลำดับที่เข้าเทียบเรียง ซึ่งวิธีนี้ใช้ได้กับการเทียบเรียงกลุ่มลำดับที่มีจำนวนน้อย และมีความยาวไม่มากนัก (วิรุฒิ, 2545) วิธีการโปรแกรมพลวัตจึงไม่เหมาะสมกับการเทียบเรียงกลุ่มลำดับข้อมูลในระดับปริมาณมากกว่า 5 ลำดับที่มีความยาวของลำดับ 100 อักขระขึ้นไป (Gusfield, 1999)

วิธีการในกลุ่มฮิวริสติกได้รับความนิยมมาก โดยเฉพาะโปรแกรม ClustalW เนื่องจากสามารถเทียบเรียงกลุ่มลำดับข้อมูลที่มีความยาว และจำนวนมากกว่าการโปรแกรมแบบพลวัต แต่หากเทียบเรียงกลุ่มลำดับข้อมูลที่มีขนาดใหญ่มากก็จะใช้เวลานานในการประมวลผลเช่นเดียวกัน

ดังนั้นการพัฒนาเทคนิคที่มีความเร็วในการประมวลผลเพิ่มขึ้นจึงเป็นสิ่งจำเป็นสำหรับปัจจุบันที่
ข้อมูลทางชีวภาพเพิ่มขึ้นอย่างรวดเร็ว

อุปกรณ์และวิธีการ

อุปกรณ์

1. ซอฟต์แวร์

- 1.1 ระบบปฏิบัติการลินุกซ์
- 1.2 โปรแกรม ClustalW เวอร์ชัน 1.83 (สามารถหาได้จาก <ftp://ftp-igbmc.u-strasbg.fr/pub/ClustalX/>)
- 1.3 โปรแกรม Bali_score สำหรับใช้วัดค่าคะแนนชุดทดสอบ BALiBASE (สามารถหาได้จาก <http://www-igbmc.u-strasbg.fr/BioInfo/BALiBASE2/>)
- 1.4 ฐานข้อมูลชีวภาพ BALiBASE เวอร์ชัน 2 (สามารถหาได้จาก <http://www-igbmc.u-strasbg.fr/BioInfo/BALiBASE2/>)
- 1.5 โปรแกรม qscore สำหรับใช้วัดค่าคะแนนชุดทดสอบ PREFAB (สามารถหาได้จาก <http://www.drive5.com/muscle/prefab.htm>)
- 1.6 ฐานข้อมูลชีวภาพ PREFAB เวอร์ชัน 4.0 (สามารถหาได้จาก <http://www.drive5.com/muscle/prefab.htm>)
- 1.7 ตัวแปลภาษา C

2. ฮาร์ดแวร์

- 2.1 เครื่องคอมพิวเตอร์ Athlon XP 2000+ ที่มี CPU Clock rate 1.67 GHz
- 2.2 หน่วยความจำหลัก 512 MB
- 2.3 ฮาร์ดดิสก์ชนิด IDE ขนาด 80 GB จำนวน 1 ตัว

วิธีการ

ในวิทยานิพนธ์นำเสนอขั้นตอนเทคนิคการเทียบเรียงกลุ่มลำดับข้อมูลโดยพิจารณาความถี่ส่วนย่อยของลำดับ (Sub-segment Frequency Alignment :SSFA) (คมสัน, 2548) ซึ่งได้แนวคิดมาจากความรู้เทคนิคการเทียบเรียงกลุ่มลำดับข้อมูลโดยใช้โปรแกรม ClustalW (Thompson *et al.*, 1999)

1. ขั้นตอนวิธีการเทียบเรียงกลุ่มลำดับข้อมูล

ขั้นตอนการเทียบเรียงกลุ่มลำดับข้อมูลโดยใช้ความถี่แบบพิจารณาค่าน้ำหนักแบ่งเป็น 3 ขั้นตอนดังภาพที่ 14



ภาพที่ 14 การเทียบเรียงกลุ่มลำดับข้อมูลโดยพิจารณาความถี่ส่วนย่อยของลำดับ

ขั้นตอนแรกเป็นการสร้างตารางค่าความแตกต่างของแต่ละคู่ลำดับข้อมูลโดยพิจารณาความถี่ส่วนย่อยซึ่งแตกต่างจากวิธีของโปรแกรม ClustalW มีรายละเอียดขั้นตอนการคำนวณดังนี้

ขั้นแรกเริ่มจากนำลำดับข้อมูลแต่ละลำดับมาแบ่งเป็นส่วนย่อยโดยใช้วิธี n-gram โดยส่วนย่อยที่ใช้ในงานวิจัยนี้มีความยาว 3 ตัวอักษร ดังตัวอย่างที่ 2

ตัวอย่างที่ 2 การแบ่งส่วนย่อยของแต่ละลำดับ

สมมติให้ข้อมูลที่ต้องการเทียบเรียงมีทั้งหมดจำนวน 3 ลำดับ ได้แก่

ลำดับ a = ACATG

ลำดับ b = TACAG

ลำดับ c = AGCATG

แบ่งเป็นส่วนย่อยได้ดังนี้

ส่วนย่อยลำดับ a คือ $a1=\{ACA\}$, $a2=\{CAT\}$, $a3=\{ATG\}$

ส่วนย่อยลำดับ b คือ $b1=\{TAC\}$, $b2=\{ACA\}$, $b3=\{CAG\}$

ส่วนย่อยลำดับ c คือ $c1=\{AGC\}$, $c2=\{GCA\}$, $c3=\{CAT\}$, $c4=\{ATG\}$

เมื่อได้ส่วนย่อยของลำดับแล้วจึงนำแต่ละส่วนย่อยของลำดับมาแปลงให้เป็นค่าตัวเลขโดยแทนค่าตัวอักษรในแต่ละหลักด้วยเลข 0-20 ตามตารางที่ 2 และนำมาผ่านสมการแฮช ดังสมการที่ 6

ตารางที่ 2 ตารางค่าตัวเลขสำหรับแทนค่าอักขระในส่วนย่อยของลำดับ

-	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

$$H(x) = (x_2 \times 21^2) + (x_1 \times 21) + x_0 \quad (6)$$

โดยที่ $H(x)$ คือค่าตัวเลขที่ได้จากการแปลงส่วนย่อยของลำดับ x ส่วน x_2 คือค่าตัวเลขที่แทนค่าอักขระตัวแรกของส่วนย่อย x , x_1 คือค่าตัวเลขที่แทนค่าอักขระตัวที่ 2 ของส่วนย่อย x และ x_0 คือค่าตัวเลขที่แทนค่าอักขระตัวที่ 3 ของส่วนย่อย x ซึ่งค่า $H(x)$ ที่คำนวณได้สำหรับแต่ละส่วนย่อยของลำดับนั้นจะได้ค่าเท่ากันก็ต่อเมื่อส่วนย่อยทั้ง 2 มีอักขระเดียวกัน และเรียงตัวเหมือนกันเท่านั้น ตัวอย่างการแปลงส่วนย่อยของลำดับเป็นดังตัวอย่างที่ 3

ตัวอย่างที่ 3 การแทนค่าส่วนย่อยของลำดับด้วยตัวเลข จากตัวอย่างที่ 2 เมื่อแบ่งลำดับออกเป็น ส่วนย่อยแล้ว นำแต่ละส่วนย่อยมาผ่านสมการแฮชได้ค่าดังนี้

ส่วนย่อยของลำดับ a ได้แก่

$$H(a1) = H(\{ACA\}) = (1 \times 21^2) + (5 \times 21) + 1 = 547$$

$$H(a2) = H(\{CAT\}) = (5 \times 21^2) + (1 \times 21) + 17 = 2243$$

$$H(a3) = H(\{ATG\}) = (1 \times 21^2) + (17 \times 21) + 8 = 806$$

ส่วนย่อยของลำดับ b ได้แก่

$$H(b1) = H(\{TAC\}) = (17 \times 21^2) + (1 \times 21) + 5 = 7523$$

$$H(b2) = H(\{ACA\}) = (1 \times 21^2) + (5 \times 21) + 1 = 547$$

$$H(b3) = H(\{CAG\}) = (5 \times 21^2) + (1 \times 21) + 8 = 2234$$

ส่วนย่อยของลำดับ c ได้แก่

$$H(c1) = H(\{AGC\}) = (1 \times 21^2) + (8 \times 21) + 5 = 614$$

$$H(c2) = H(\{GCA\}) = (8 \times 21^2) + (5 \times 21) + 1 = 3634$$

$$H(c3) = H(\{CAT\}) = (5 \times 21^2) + (1 \times 21) + 17 = 2243$$

$$H(c4) = H(\{ATG\}) = (1 \times 21^2) + (17 \times 21) + 8 = 806$$

เมื่อแปลงส่วนย่อยของลำดับเป็นตัวเลขโดยผ่านฟังก์ชันแฮชแล้วจึงเก็บข้อมูลลงในตารางแบบแฮช ส่วนย่อยที่เหมือนกันจะถูกเก็บอยู่ในตำแหน่งของตารางแฮชเดียวกัน โดยข้อมูลที่เก็บในตารางแฮชจะมีการเก็บจำนวนของส่วนย่อยของแต่ละลำดับข้อมูลไว้ด้วย

ในระหว่างการเก็บข้อมูลส่วนย่อยลงในตารางแฮชนั้น จะมีการเก็บข้อมูลความถี่ของส่วนย่อยที่ซ้ำกันของแต่ละลำดับข้อมูลเพื่อสร้างตารางค่าความคล้ายกันของแต่ละลำดับข้อมูล ซึ่งค่าความคล้ายกันคำนวณจากอัตราส่วนของความถี่ของส่วนย่อยที่ซ้ำกันของลำดับข้อมูลต่อจำนวนส่วนย่อยทั้งหมดของลำดับข้อมูล ดังสมการที่ 7

$$S_{SSFA}(x, y) = \frac{\sum_{\mathcal{L}} (k(\mathcal{L}))}{(n_x + n_y)} \quad (7)$$

โดยที่ $S_{SSFA}(x, y)$ คือค่าคะแนนความคล้ายกันพิจารณาจากความถี่ของส่วนย่อยที่ซ้ำกันของลำดับ x และ y , \mathcal{L} คือส่วนย่อย n -gram ที่ซ้ำกัน, n_x และ n_y คือจำนวนของส่วนย่อยแบบ 3-gram ของลำดับ x และ y และ $k(\mathcal{L})$ คือความถี่ของส่วนย่อย \mathcal{L} ที่ปรากฏในลำดับ x และ y

ค่าความคล้ายกันของคู่ลำดับข้อมูลมีค่าอยู่ระหว่าง 0 ถึง 1 โดยค่า 0 หมายถึงคู่ลำดับข้อมูลที่ไม่มีความคล้ายกันเลย และ 1 คือคู่ลำดับข้อมูลที่เหมือนกันทุกประการ ตัวอย่างการคำนวณค่าความคล้ายกันของคู่ลำดับข้อมูลเป็นดังตัวอย่างที่ 4

ตัวอย่างที่ 4 การคำนวณค่าความคล้ายกันสำหรับแต่ละคู่ลำดับข้อมูลจากตัวอย่างที่ 3 สามารถคำนวณได้โดยใช้สมการที่ 7 ดังนี้

ค่าความคล้ายกันของลำดับ a และลำดับ b คือ

$$S_{SSFA}(a,b) = \frac{k(\{ACA\})}{n_a + n_b} = \frac{2}{3+3} = 0.333$$

ค่าความคล้ายกันของลำดับ a และลำดับ c คือ

$$S_{SSFA}(a,c) = \frac{k(\{CAT\}) + k(\{ATG\})}{n_a + n_c} = \frac{2 + 2}{3 + 4} = 0.571$$

ค่าความคล้ายกันของลำดับ b และลำดับ c คือ

$$S_{SSFA}(b,c) = \frac{0}{n_b + n_c} = \frac{0}{3 + 4} = 0$$

เมื่อได้ค่าความคล้ายกันของแต่ละคู่ลำดับข้อมูลแล้ว จึงสร้างตารางค่าความแตกต่างโดยคำนวณจากค่าความคล้ายกันดังสมการที่ 8 และจากตัวอย่างที่ 4 เมื่อคำนวณค่าความแตกต่างแล้วจะได้ตารางค่าความแตกต่างของแต่ละคู่ลำดับข้อมูลดังตารางที่ 3

$$D_{SSFA}(x,y) = 1 - S_{SSFA}(x,y) \quad (8)$$

เมื่อ $D_{SSFA}(x,y)$ คือค่าความแตกต่างของลำดับ x และลำดับ y

ตารางที่ 3 ตารางค่าความแตกต่างของลำดับ a, b และ c ที่เข้าเทียบเรียง

	ลำดับ a	ลำดับ b	ลำดับ c
ลำดับ a	0	0.667	0.429
ลำดับ b	0.667	0	1
ลำดับ c	0.429	1	0

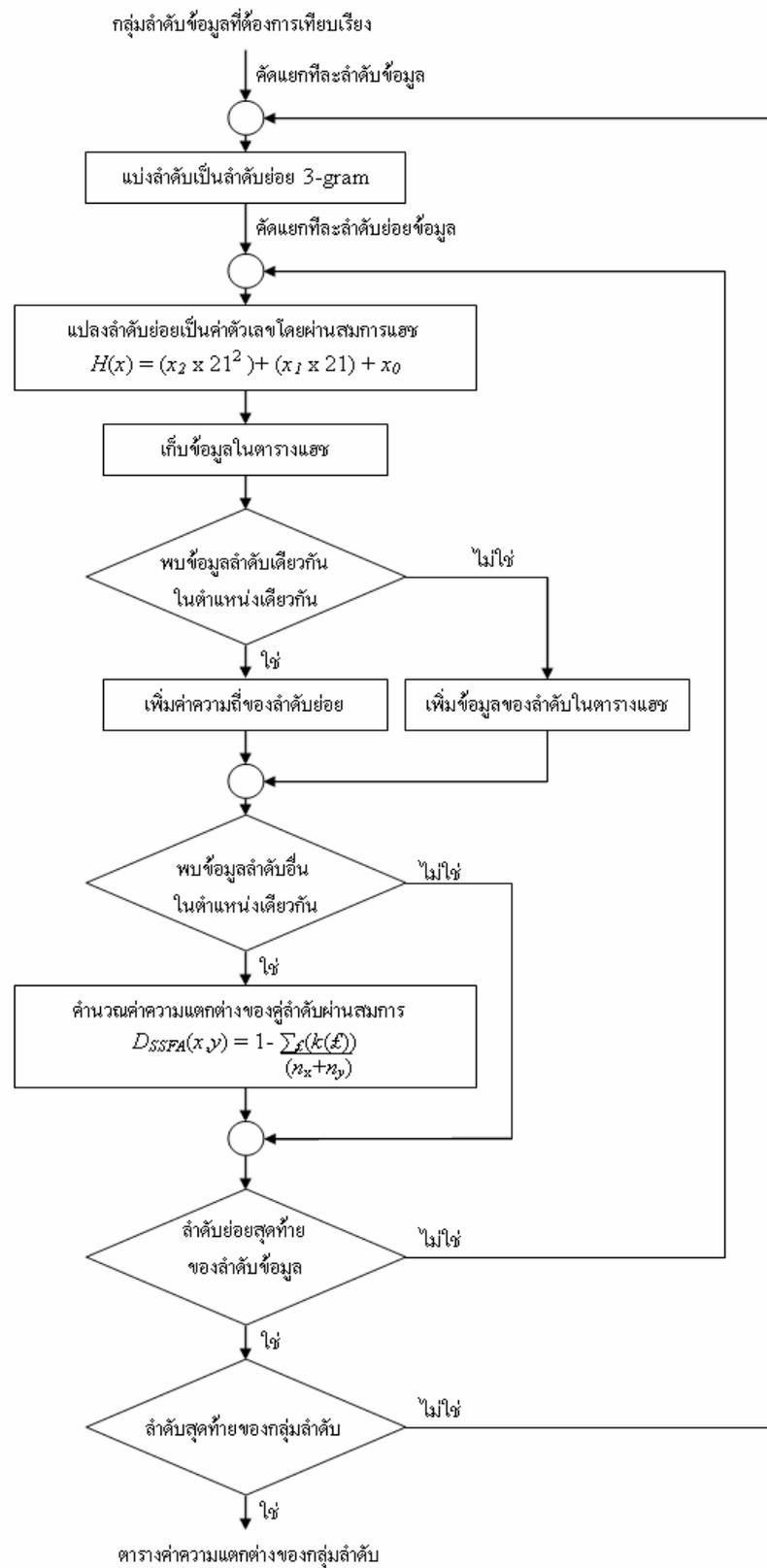
การคำนวณค่าความแตกต่างสามารถคำนวณสะสมได้พร้อมกับการเก็บข้อมูลแต่ละส่วนย่อยในตารางแฮช โดยใช้สมการที่ 9 ซึ่งเกิดจากการแทนค่า $S_{SSFA}(x,y)$ ในสมการที่ 8 ด้วยสมการที่ 7 ดังนั้นขั้นโปรแกรม SSFA จึงมีขั้นตอนการทำงานดังภาพที่ 15 และมีตัวอย่างการคำนวณดังภาพที่ 16

$$D_{SSFA}(x, y) = 1 - \frac{\sum_{\ell} (k(\ell))}{(n_x + n_y)} \quad (9)$$

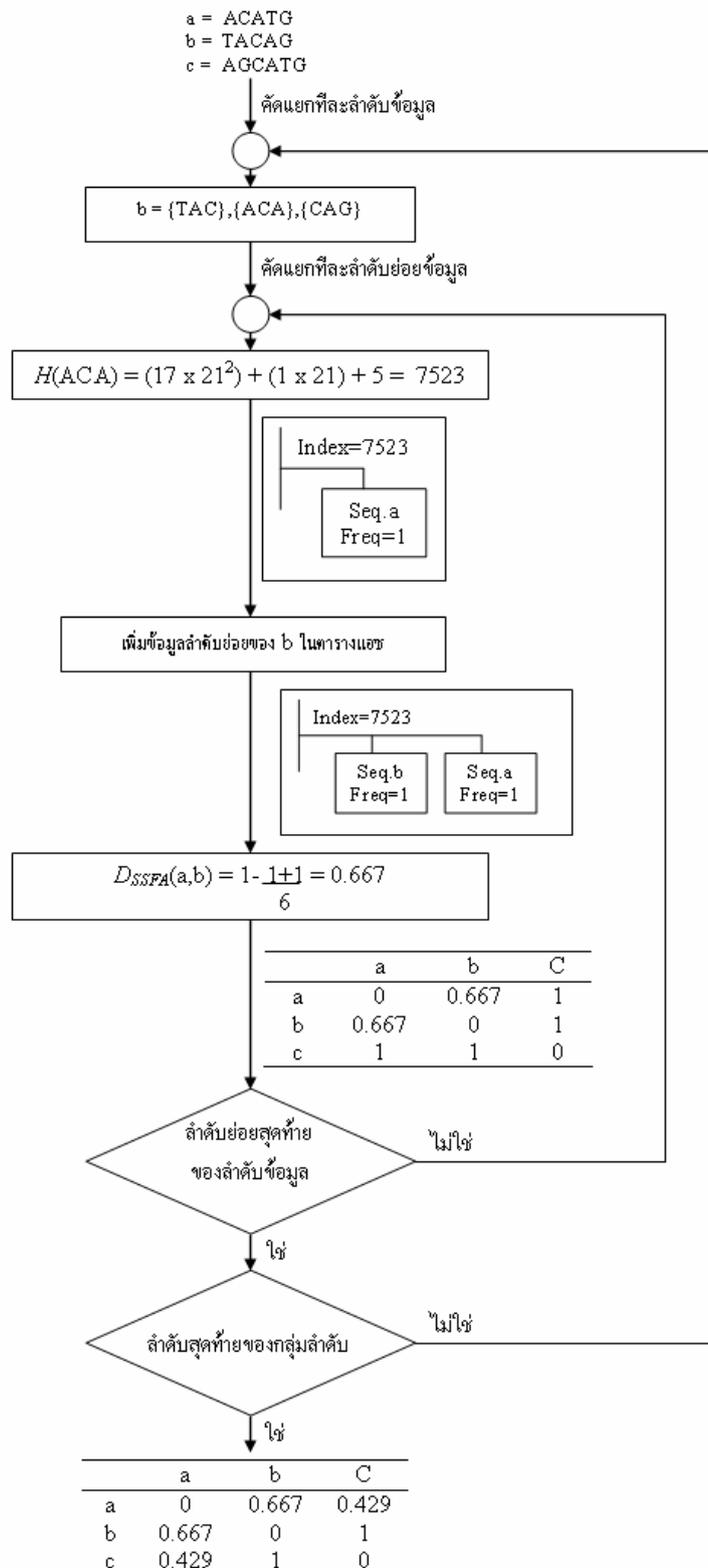
เมื่อได้ตารางค่าความแตกต่างของแต่ละคู่ลำดับข้อมูลแล้ว จึงทำขั้นตอนจัดกลุ่มและลำดับที่เข้าเทียบเรียง และขั้นตอนเทียบเรียงกลุ่มลำดับข้อมูลแบบก้าวหน้าดังภาพที่ 14 โดยทั้ง 2 ขั้นตอนนี้ใช้วิธีการเดียวกับโปรแกรม ClustalW แต่อ้างอิงจากตารางค่าความแตกต่างที่คำนวณจากความถี่ส่วนย่อยของข้อมูลแทนการเทียบเรียงคู่ลำดับของโปรแกรม ClustalW

การแบ่งส่วนย่อยลำดับข้อมูลจำนวน N ลำดับ โดยวิธี n-gram นั้นมีค่าความซับซ้อนของเวลา $O(N(L-n+1))$ หรือประมาณ $O(NL)$ เมื่อ L คือความยาวของลำดับข้อมูล และ n คือความยาวส่วนย่อย n-gram ในส่วนของการคำนวณค่าความคล้ายกันนั้นใช้วิธีคำนวณแบบสะสมจากการเก็บข้อมูลใน hash เดียว ทำให้การสร้างตารางค่าความแตกต่างของคู่ลำดับข้อมูลนั้นมีการประมวลผลแต่ละลำดับข้อมูลเพียงรอบเดียว

การทำงานของโปรแกรม SSFA และการทำงานของโปรแกรม ClustalW นั้นใช้เวลาในการประมวลผลแตกต่างกันในขั้นตอนการสร้างตารางค่าความแตกต่างของคู่ลำดับข้อมูล โดยโปรแกรม SSFA ใช้วิธีการพิจารณาความถี่ของส่วนย่อย และมีค่าซับซ้อนของเวลา $O(NL)$ ในขณะที่โปรแกรม ClustalW ใช้วิธีเทียบเรียงแต่ละคู่ลำดับข้อมูล และมีค่าซับซ้อนของเวลา $O(N^2L^2)$ ซึ่งทำให้โปรแกรม SSFA ใช้เวลาเทียบเรียงกลุ่มลำดับข้อมูลน้อยกว่าโปรแกรม ClustalW มาก



ภาพที่ 15 ขั้นตอนการสร้างตารางค่าความแตกต่างโดยพิจารณาความถี่ส่วนย่อย



ภาพที่ 16 ขั้นตอนการสร้างตารางค่าความแตกต่างโดยพิจารณาความถี่ส่วนย่อย

2. การเตรียมข้อมูล

ฐานข้อมูลที่ใช้ในงานวิจัยนี้ คือฐานข้อมูล BALiBASE เวอร์ชัน 2.01 และฐานข้อมูล PREFAB เวอร์ชัน 4.0

ฐานข้อมูล BALiBASE เวอร์ชัน 2.01 (Thompson *et al.*, 1999; วิรุฒิ, 2545) เป็นฐานข้อมูลสำหรับวัดผลการเทียบเรียงลำดับข้อมูล โปรตีน ประกอบด้วยกลุ่มลำดับข้อมูลที่เทียบเรียงแล้วสำหรับอ้างอิง 141 กลุ่ม แบ่งเป็นกลุ่มอ้างอิง 8 ชุด โดยในงานวิจัยนี้ใช้กลุ่มอ้างอิงทั้งหมด 5 กลุ่มที่สามารถวัดผลการเทียบเรียงได้ คือ ชุดอ้างอิงที่ 1 เป็นชุดที่มีค่าความห่างเท่ากันและหลายระดับส่วน, ชุดอ้างอิงที่ 2 เป็นการเทียบเรียงกับสายอักขระที่มีการเปลี่ยนแปลงสูง, ชุดอ้างอิงที่ 3 เป็นกลุ่มย่อยที่น้อยกว่า 25% ของส่วนที่เหลือระหว่างกลุ่ม, ชุดอ้างอิงที่ 4 แต่ละกลุ่มข้อมูลประกอบด้วยลำดับข้อมูลที่เป็นส่วนขยายขนาดใหญ่ ที่จุดสิ้นสุด N/C(N/C-terminal extensions) และชุดอ้างอิงที่ 5 แต่ละกลุ่มข้อมูลประกอบด้วยลำดับข้อมูลที่เกิดการแทรกภายในช่วงกลางลำดับ โดยมีรายละเอียดแต่ละชุดอ้างอิงดังตารางที่ 4 – ตารางที่ 9

ตารางที่ 4 ชุดอ้างอิง 1 เป็นชุดที่มีค่าความห่างเท่ากันและหลายระดับส่วนสวาง (ขนาดสั้น)

คุณลักษณะ	ชื่อ(จำนวนสายอักขระ)
ความเหมือนกัน(identity) น้อยกว่า 25%	laboA (5), lidy (5), lr69 (4), ltvxA (4), lubi (4), lwit (5), 2trx (4) รวม 31 สายอักขระ
ความเหมือนกัน 20-40%	laab (4), lfj1A (6), lhfh (5), lhpi (4), lcsy (5), lpf (5), ltgx (4), lycc (4), 3cyr (4), 451c (5) รวม 46 สายอักขระ
ความเหมือนกันมากกว่า 35%	laho (5), lisp (5), ldox (4), lfkj (5), lfmb (4), lkm (5), lplc (5), 2fxb (5), 2mhr (5), 9mt (5) รวม 48 สายอักขระ

ตารางที่ 5 ชุดอ้างอิง 1 เป็นชุดที่มีค่าความห่างเท่ากันและหลายระดับส่วนสวาง(ขนาดกลาง)

คุณลักษณะ	ชื่อ(จำนวนสายอักขระ)
ความเหมือนกัน(identity) น้อยกว่า 25%	lbbt3 (5), lsbp (5), lhavA (5), luky (4), 2hsdA (4), 2pia (4), 3grs (4), kinase (5) รวม 36 สายอักขระ
ความเหมือนกัน 20-40%	lad2(4), laym3(4), lgdoA(4), lldg(4), lmrj(4), lpgtA(4), lpii(4) รวม 38 สายอักขระ
ความเหมือนกันมากกว่า 35%	lamk (5), lar5A (4), lezm (5), lled (4), lppn (5), lpysA (4), lthm (4), ltis (5), lzin (4), 5ptp (5) รวม 45 สายอักขระ

ตารางที่ 6 ชุดอ้างอิง 1 เป็นชุดที่ค่าความห่างเท่ากันและหลายระดับส่วนสงวน (ขนาดยาว)

คุณลักษณะ	ชื่อ(จำนวนสายอักขระ)
ความเหมือนกัน (identity) น้อยกว่า 25%	lajsA (4), 1cpt (4), 1lvl (4), 1pamA (5), 1ped (3), 2myr (4), 4enl (3), gal4 (5) รวม 32 สายอักขระ
ความเหมือนกัน 20-40%	1ac5 (4), 1adj (4), 1bgl (4), 1dlc (4), 1eft (4), 1fieA (4), 1gowA (4), 1pkm (4), 1sesA (5), 2ack (5), arp (5), glg (5) รวม 52 สายอักขระ
ความเหมือนกันมากกว่า 35%	1ad3 (4), 1gpb (5), 1gr (5), 1lcf (6), 1rthA (5), 1taq (5), 3pmg (4), actin (5) รวม 39 สายอักขระ

ตารางที่ 7 ชุดอ้างอิง 2 เป็นการเทียบเรียงกับสายอักขระที่มีการเปลี่ยนแปลงสูง

คุณลักษณะ	ชื่อ(จำนวนสายอักขระ)
ขนาดสั้น	1aboA (15), 1csy (19), 1idy (21), 1r69 (22), 1tgxA (20), 1tvxA (17), 1ubi (17), 1wit (22), 2trx (20) รวม 173 สายอักขระ
ขนาดกลาง	1sbp (18), 1havA (18), 1uky (24), 2hsdA (22), 2pia (18), 3grs (16), kinase (18) รวม 134 สายอักขระ
ขนาดยาว	lajsA (20), 1cpt (15), 1lvl (24), 1pamA (19), 1ped (19), 2myr (19), 4enl (18) รวม 134 สายอักขระ

ตารางที่ 8 ชุดอ้างอิง 3 เป็นกลุ่มย่อยที่น้อยกว่า 25% ของส่วนที่เหลือระหว่างกลุ่ม

คุณลักษณะ	ชื่อ(จำนวนสายอักขระ)
ขนาดสั้น	1idy (27), 1r69 (23), 1ubi (22), 1wit (19) รวม 91 สายอักขระ
ขนาดกลาง	1uky (24), 2pia (20), kinase (23) รวม 67 สายอักขระ
ขนาดยาว	lajsA (28), 1pamA (19), 1ped (21), 2myr (21), 4enl (19) รวม 108 สายอักขระ

ตารางที่ 9 ชุดอ้างอิง 4 และชุดอ้างอิง 5

ชุดอ้างอิง	คุณลักษณะ	ชื่อ(จำนวนสายอักขระ)
4	สายอักขระ N/C (N/C-terminal extensions)	lckaA (10), lcsp (6), ldynA (16), llkl (8), lmfa (8), lpfC (10), lpysA (4), lvln (14), lycc (9), 2abk (7), kinase1 (7), kinase2 (18) รวม 107 สายอักขระ
5	สายอักขระที่มีการ แทรกภายใน	left (8), livy (7), lpysA (10), lppg (5), lthm1 (11), lthm2 (7), 2cba (8), s51 (15), s52 (5), kinase1 (5), kinase2 (12), kinase3 (19) รวม 112 สายอักขระ

ข้อมูลทดสอบมาจาก <http://www-igbmc.u-strasbg.fr/BioInfo/BALiBASE2/>

BALiBASE2.01.tar.gz มีขนาด 7.81 MB และเมื่อทำการขยายไฟล์จะเป็นชุดทดสอบที่มีขนาด 82.4 MB ประกอบด้วย 8 กลุ่มย่อยตามชุดอ้างอิง โดยมีชื่อ ref1, ref2, ref3, ref4, ref5, ref6, ref7 และ ref8 ตามลำดับ ในงานวิจัยนี้ใช้กลุ่มย่อย ref1 – ref5 โดยใช้ชุดอ้างอิงเท่ากับ ref เช่น ชุดอ้างอิง: 1 ตามโครงสร้างคือ ref1 และชุดทดสอบเท่ากับ test เช่น ชุดทดสอบ 1 ตามโครงสร้างคือ test1 ชื่อชุดข้อมูลแสดงได้ดังตารางที่ 10 และตารางที่ 11

ตารางที่ 10 ชุดอ้างอิง 1

ชุดทดสอบ	ชื่อชุดข้อมูล
1	laab_ref1, laboA_ref1, laho_ref1, lcsp_ref1, lcsy_ref1, ldox_ref1, lfj1A_ref1, lfkj_ref1, lfmb_ref1, lhfh_ref1, lhpi_ref1, lidy_ref1, lkrm_ref1, lpfC_ref1, lplc_ref1, lr69_ref1, ltgxA_ref1, ltvxA_ref1, lubi_ref1, lwit_ref1, lycc_ref1, 2fxb_ref1, 2mhr_ref1, 2trx_ref1, 3cyr_ref1, 451c_ref1, 9rnt_ref1
2	lad2_ref1, lamk_ref1, lar5A_ref1, laym3_ref1, lbbt3_ref1, lezm_ref1, lgdoA_ref1, lhavA_ref1, lldg_ref1, lled_ref1, lmrj_ref1, lpgt_ref1, lpii_ref1, lppn_ref1, lpysA_ref1, lsbp_ref1, lthm_ref1, ltis_ref1, lton_ref1, luky_ref1, lzin_ref1, 2cba_ref1, 2hsdA_ref1, 2pia_ref1, 3grs_ref1
3	lac5_ref1, lad3_ref1, ladj_ref1, lajsA_ref1, lbgl_ref1, lcpt_ref1, ldc_ref1, left_ref1, lfieA_ref1, lgowA_ref1, lgpB_ref1, lgtr_ref1, llcf_ref1, llvl_ref1, lpamA_ref1, lped_ref1, lpkm_ref1, lrthA_ref1, lsesA_ref1, ltaq_ref1, 2ack_ref1, 2myr_ref1, 3pmg_ref1, 4enl_ref1, actin_ref1, arp_ref1, gal4_ref1, glg_ref1

ตารางที่ 11 ชุดอ้างอิง 2 - 5

ชุดอ้างอิง	ชื่อชุดข้อมูล
2	laboA_ref2, lajsA_ref2, lcpt_ref2, lcsy_ref2, lhavA_ref2, lidy, lidy_ref2, llvl_ref2, lpamA_ref2, lped_ref2, lr69_ref2, lsbp_ref2, ltgxA_ref2, ltvxA_ref2, lubi_ref2, luky_ref2, lwit_ref2, 2hsdA_ref2, 2myr_ref2, 2pia_ref2, 2trx_ref2, 3grs_ref2, 4enl_ref2, kinase_ref2
3	lajsA_ref3, lidy_ref3, lpamA_ref3, lped_ref3, lr69_ref3, lubi_ref3, luky_ref3, lwit_ref3, 2myr_ref3, 2pia_ref3, 4enl_ref3, kinase_ref3
4	lckaA_ref4, lcsp_ref4, ldynA_ref4, llkl_ref4, lmfa_ref4, lpfc_ref4, lpysA_ref4, lvln_ref4, lycc_ref4, 2abk_ref4, kinase1_ref4, kinase2_ref4
5	left_ref5, livy_ref5, lpysA_ref5, lqpg_ref5, lthm1_ref5, lthm2_ref5, 2cba_ref5, kinase1_ref5, kinase2_ref5, kinase3_ref5, s51_ref5, s52_ref5

ฐานข้อมูล PREFAB (Edgar, 2004) เวอร์ชัน 4.0 ประกอบด้วยกลุ่มข้อมูลสำหรับเทียบเรียง 1682 ชุดข้อมูล โดยไม่ได้แบ่งเป็นกลุ่มอ้างอิงเช่นเดียวกับ BAliBASE แต่ละชุดข้อมูล ประกอบด้วยลำดับข้อมูลโปรตีนที่มาจากการค้นหาลำดับข้อมูลที่เทียบเคียงกันได้ผ่านโปรแกรม PSI-BLAST จากฐานข้อมูลโปรตีน NCBI จำนวนลำดับข้อมูลในแต่ละชุดข้อมูลโดยเฉลี่ยมีจำนวน 45 ลำดับ

ฐานข้อมูล PREFAB เวอร์ชัน 4.0 มาจาก <http://www.drive5.com/muscle/prefab.htm> โดยมีขนาด 7.6 MB เมื่อขยายไฟล์ออกมาจะเป็นชุดทดสอบที่มีขนาด 31 MB ประกอบด้วยชุดข้อมูล 1682 ชุด ในงานวิจัยนี้จะแบ่งออกเป็นกลุ่มข้อมูลจำนวน 68 กลุ่มดังตารางที่ 12 – ตารางที่ 21

ตารางที่ 12 กลุ่มข้อมูลที่ 1-7

กลุ่มข้อมูล	ชื่อชุดข้อมูล
1	12asA_1atiA, 1531_1cnsA, 16pk_1qpg, 19hcA_1duwA, 1a02N_1a3qA, 1a04A_1dz3A, 1a04A_1ibjA, 1a04A_3chy, 1a0cA_1a0dA, 1a0cA_1bxbd, 1a0cA_1ubpC, 1a0fA_1hqoA, 1a0fA_1ljrA, 1a0fA_1pd21, 1a0hA_5hpgA, 1a0p_1a36A, 1a0p_1ae9A, 1a0tP_2mprA, 1a12A_1jtdB, 1a1z_1d2zA, 1a28A_2prgA, 1a28A_3erdA, 1a2vA_1ksiB, 1a36A_1a41, 1a3aA_1a6jA
2	1a3c_1tc1A, 1a3k_1c11A, 1a3k_1lcl, 1a3qA_1gof, 1a44_1qouB, 1a49A_1dxeA, 1a49A_1pkyC, 1a49A_2tpsA, 1a4iA_1b0aA, 1a4iA_1ee9A, 1a53_1dvjA, 1a53_1hg3A, 1a53_1nsj, 1a53_1qo2A, 1a62_1mjc, 1a65A_1aozA, 1a65A_1nif, 1a6cA_1bmv1, 1a6dA_1ass, 1a6dA_1derA, 1a6jA_1hynP, 1a6l_2fdn, 1a6m_1ash, 1a6m_1cg5B, 1a6m_1ewaA
3	1a6m_1flp, 1a6m_1h97A, 1a6m_1h1b, 1a6m_1lithA, 1a6m_2fal, 1a6m_2gdm, 1a6m_2hbg, 1a6m_2vhhA, 1a6m_3sdhA, 1a6o_1buhA, 1a75A_1b8cB, 1a7s_1agjA, 1a7s_2cgaB, 1a7tA_1e5dA, 1a7tA_1smlA, 1a7w_1aoiA, 1a7w_1aoiB, 1a7w_1b67B, 1a7w_1bh9B, 1a7w_1tafA, 1a7w_1tafB, 1a8d_1wba, 1a8h_1qqtA, 1a8l_1hyuA, 1a8o_1qrjB
4	1a8rA_1fb1D, 1a9nA_1dceA, 1aba_1h75A, 1aba_1kte, 1aba_1qfnA, 1aba_3grx, 1abrB_1dqgA, 1abwA_1cqxA, 1ac5_1auoA, 1ac5_1cpy, 1ad3A_1bpwA, 1ad3A_1euhA, 1adeA_1dj3A, 1adjA_1atiA, 1adjA_1httD, 1ae9A_1a36A, 1ae9A_1aihA, 1aerA_1aerB, 1ag4_1amm, 1agdA_1bqsA, 1agdA_1iakA, 1agdB_1iakA, 1agjA_1cqqa, 1agqA_1vpfA, 1ahl_1cid
5	1ahl_1f97A, 1ahl_1kb5B, 1ahl_1qfoA, 1ahl_1tlk, 1ahl_1vcaA, 1ahl_2ncm, 1ah7_1ca1, 1ah9_1kbb, 1aihA_1a0p, 1air_1pcl, 1aisB_1guxB, 1aj6_1b63A, 1aj8A_1csh, 1ajqA_1cp9A, 1ajqB_1nedA, 1ajsA_1aam, 1ajsA_1ibjA, 1ak1_1hrkA, 1ako_1bix, 1alu_1bgc, 1alu_1cnt3, 1alu_1i1rB, 1aly_1d4vB, 1aly_1dyoA, 1aly_2tnfA
6	1am2_1at0, 1am7A_3lzt, 1amj_1zymA, 1amk_1hg3A, 1amoA_1ddgB, 1amoA_1ep2B, 1amoA_1fnc, 1amoA_1qfjA, 1amoA_1rcf, 1amp_1b8oA, 1amp_1cg2A, 1amp_1cx8B, 1amp_1xjo, 1amuA_1lci, 1an8_1aw7A, 1an8_1eu4A, 1an8_1f77B, 1an9A_1c0pA, 1aoa_1bkrA, 1aocA_1bndA, 1aoeA_1d1gA, 1aoeA_1vdrA, 1aoeA_3dfr, 1aohA_1anu, 1aohA_1g1kB
7	1aoiA_1bh9B, 1aoiA_1tafA, 1aoiB_1a7w, 1aoiB_1aoiD, 1aoiB_1jfiB, 1aoiB_1tafA, 1aoiC_1f66G, 1aoiD_1aoiB, 1aoiD_1jfiB, 1aoxA_1atzA, 1aoxA_1auq, 1aoxA_1ido, 1aoy_1b4aA, 1ap0_1dz1A, 1ap8_1ejhA, 1apyA_1ayyC, 1apyB_2gawD, 1aq0A_1eokA, 1aq0A_1ghsB, 1aqb_1avgI, 1aqb_1bebA, 1aqb_1epaA, 1aqb_1h91A, 1aqt_1e79H, 1aquA_1fmlB

ตารางที่ 13 กลุ่มข้อมูลที่ 8-14

กลุ่มข้อมูล	ชื่อชุดข้อมูล
8	laqzA_1de3A, laqzA_9rnt, larb_2hrvA, lash_1a6m, lash_1flp, lash_1lithA, lash_2fal, lat3A_1fl1A, lat3A_1lay, latg_1ixh, latg_1mrp, latg_1wod, latiA_1g5hA, latiA_1qf6A, latlA_1bkcA, latzA_1auq, latzA_1ido, latzA_2scuB, lau7A_2hddA, lauiB_1dguA, lauk_1fsu, lauoA_1brt, lauoA_1fj2B, lauoA_1i6wA, lauoA_1jfrA
9	lauq_1ido, lauvA_1iow, lauyA_1e57B, lauyA_1qqp3, lauz_1dciA, lavaC_1wba, lavgI_1epaA, lavmA_3mdsA, lavpA_1euvA, law0_1cc8A, law8B_1cr5A, lawcA_1bc8C, lawcA_2hts, lawcA_2irfG, lawd_1ayfA, lawe_1btkA, lawe_1dbhA, lawe_1faoA, lawe_1pls, lax4A_2dkb, lax4A_2tplA, lax8_1f6fA, lax8_1rcb, laxiB_1bp3B, laxiB_1j7vR
10	laxkA_2nlrA, laye_1qmuA, layfA_1put, laym1_2hfw1, laym2_1hri2, laym2_1qqp2, laym3_1b35A, laym3_1ihmA, laym3_1rud3, laz9_1bn5, laz9_1c24A, lazpA_1bf4A, lazsa_1culB, lazsa_1fx2A, lb0nA_1r69, lb0uA_1jj7A, lb0uA_1skyE, lb16A_1bdb, lb16A_1bsvA, lb16A_1e6wA, lb16A_1fds, lb16A_1he2A, lb16A_1xel, lb16A_1ybvA, lb16A_3chy
11	lb20A_1rgeA, lb25A_1aorA, lb2pA_1msaA, lb34B_1d3bA, lb34B_1d3bB, lb34B_1i8fA, lb35A_1b35B, lb35A_1ihmA, lb35A_2mev1, lb37B_1gpeA, lb3aA_1hfgA, lb3aA_1tvxA, lb3jA_1agdA, lb3mA_1b8sA, lb4kA_1qmlA, lb5fB_1smrA, lb5l_1rmi, lb64_1fjFJ, lb6cB_1ia8A, lb6e_1fm5A, lb6e_1hyrB, lb6g_1bn6A, lb6g_1cqWA, lb6rA_1gsoA, lb6tA_1f9aA
12	lb7eA_1f3iA, lb7fA_1b64, lb7gO_1drw, lb7gO_1gcuA, lb7gO_1ofgA, lb87A_1i21A, lb8aA_1aszA, lb8aA_1lylA, lb8gA_1bjwA, lb8gA_1bw0A, lb8gA_1iaxA, lb8oA_1cb0A, lb8xA_1pd21, lb8xA_1pgtA, lb93A_1g8mA, lb9hA_1ibjA, lb9IA_1dhn, lb9yC_1a0rP, lb9zA_1a49A, lba1_1dkgD, lba1_1glcG, lbag_1smd, lbak_1btkA, lbak_1btn, lbak_1dynA
13	lbak_1faoA, lbak_1mai, lbb9_2semA, lbbhA_1jafB, lbbhA_2ccyA, lbbzA_1ckaA, lbbzA_1gcpA, lbcfA_1dpsA, lbcfA_1eumA, lbcfA_1qghA, lbcfA_2fha, lbcpA_1lt3A, lbcpD_1prtF, lbcpD_1tiiD, lbcpD_3chbD, lbd3A_1dkuA, lbd3A_1nulA, lbdb_1cydA, lbdb_1e6wA, lbdb_1e7wA, lbdb_1hdcA, lbdb_1oaa, lbdb_1ybvA, lbdo_1ghj, lbdo_1iyu
14	lbdyA_1rlw, lbe3A_1be3B, lbe3A_1ezvA, lbe3A_1ezvB, lbe3A_1hr6A, lbe3B_1bccB, lbe3B_1ezvB, lbe3B_1hr6A, lbe3C_3bccC, lbebA_1bj7, lbebA_1epaA, lbebA_1mup, lbec_1iakB, lbec_1tvdA, lbec_3mcg1, lbefA_1jxpA, lbefA_2hrvA, lbf2_1ehaA, lbfd_1d4oA, lbfG_1afcH, lbfG_1jlyA, lbg2_1f9tA, lbgIA_1bhgA, lbh9B_1jfiB, lbh9B_1tafA

ตารางที่ 14 กลุ่มข้อมูลที่ 15 – 21

กลุ่มข้อมูล	ชื่อชุดข้อมูล
15	1bhe_1czfA, 1bhe_1rmg, 1bhgA_1qnoA, 1bhtA_5hpgA, 1bi0_1fx7B, 1bi0_1smtA, 1bi5A_1dd8A, 1bi5A_1hn9A, 1bif_1qhfA, 1bj4A_1bjwA, 1bj7_1epaA, 1bj7_1obpA, 1bjnA_1bt4A, 1bjwA_1bw0A, 1bjwA_1cl2A, 1bjwA_1d2fA, 1bjx_1qgvA, 1bk5A_1ialA, 1bk7A_1bolA, 1bkjA_1f5vA, 1bkjA_1vfrA, 1bkpA_1tlcA, 1bkrA_1aoa, 1bkrA_1bhdA, 1bli_1jdc
16	1bm9A_1cf7B, 1bmdA_1b8pA, 1bmdA_1ceqA, 1bmdA_1hyhA, 1bmdA_2cmd, 1bmfG_1mabG, 1bmlC_1qqrC, 1bn5_1c24A, 1bn6A_1cv2A, 1bndA_1aocA, 1bndA_1nt3A, 1bo4A_1cjwA, 1bo4A_1i21A, 1bob_1qsmA, 1bob_1yghA, 1bolA_1bk7A, 1booA_1eg2A, 1boy_1a21A, 1bp7A_1af5, 1bpi_1brcl, 1bpv_1bquA, 1bpv_2fnbA, 1bpwA_1euhA, 1bpwA_1eyyA, 1bqcA_7a3hA
17	1bqg_1ec7B, 1bqg_1mucA, 1bqg_1qumA, 1bqg_2mnr, 1bqk_2mtaA, 1bqsA_1cvsC, 1bquA_1bj8, 1br9_1d2bA, 1br9_1jc7A, 1brmA_1drw, 1brmA_1gcuA, 1brmA_1qrrA, 1brt_1i6wA, 1brt_1jfrA, 1brt_1qlwA, 1bs0A_1b8gA, 1bs0A_1elqA, 1bs9_1cex, 1bs9_1ei9A, 1bsvA_1bwsA, 1bsvA_1bxkA, 1bsvA_1he2A, 1bsvA_1xel, 1bt3A_1lla, 1btka_1btn
18	1btka_1dynA, 1btka_1faoA, 1btka_1pls, 1btka_1rrpB, 1btl_1e25A, 1btl_1ei5A, 1btl_1omeB, 1btl_1skf, 1btn_1dynA, 1btn_1mai, 1btn_1pls, 1btn_1qqgA, 1btn_1rrpB, 1bu2A_1g3nC, 1bu2A_1vin, 1bu7A_1ea1A, 1burA_1rusB, 1bv1_1em2A, 1bv1_1pmgA, 1bvwa_1dysA, 1bvwa_1tml, 1bvza_1cxlA, 1bvza_1ehaA, 1bvza_1smaA, 1bvza_7taa
19	1bw0A_1d2fA, 1bw9A_1qorA, 1bw9A_1qp8A, 1bx4A_1dgyA, 1bx4A_1rkd, 1bxkA_1db3A, 1bxkA_1eq2A, 1bxkA_1he2A, 1bxkA_1xel, 1byb_1fa2A, 1byfA_1esl, 1byfA_1fm5A, 1byfA_1htn, 1byfA_1rtm1, 1byi_1eg7A, 1byi_1nksA, 1bykA_1gca, 1bykA_1qpzA, 1bykA_1tlfA, 1bykA_2dri, 1bylA_1ecsA, 1bylA_1qipA, 1bylA_1qtoA, 1byuA_1ctqA, 1bywA_1drmA
20	1bywA_1g28D, 1bywA_3pyp, 1c05A_1dm9A, 1c0aA_1adjA, 1c0aA_1lylA, 1c0nA_1ecxA, 1c0nA_1elqA, 1c20A_1ig6A, 1c24A_1chmA, 1c24A_1xgmA, 1c3cA_1fl0A, 1c3d_1qsjD, 1c3oB_1qdlB, 1c3wA_1brd, 1c3wA_1jgjA, 1c4rA_1qu0A, 1c4rA_1sacA, 1c4xA_1jfrA, 1c53_2dvh, 1c8bA_1cfzA, 1c8oA_1hleA, 1c8oA_1ovaA, 1c8oA_1sek, 1c8pA_1cto, 1c9kB_1cbuC
21	1c9kB_1d2nA, 1c9kB_1esc, 1c9kB_1g5rA, 1c9kB_2dhqA, 1ca1_1ah7, 1caxB_1dgwA, 1cb8A_1eguA, 1ccza_1fltX, 1ccza_1hngB, 1ccza_1qfoA, 1ccza_1tlk, 1cd31_1cd32, 1cd8_1hxmD, 1cd8_1tvdA, 1cd8_1wit, 1cdkA_1apmE, 1cdkA_1csn, 1cdkA_1hcl, 1cdkA_1ia8A, 1ceo_1cz1A, 1ceo_1edg, 1ceqA_1hyhA, 1ceqA_2cmd, 1ceqA_3ldh, 1cewI_1eqkA

ตารางที่ 15 กลุ่มข้อมูลที่ 22 - 28

กลุ่มข้อมูล	ชื่อชุดข้อมูล
22	lcf7A_lcf7B, lcf9A_lg2iA, lcg2A_lcx8B, lcg2A_lxjo, lcg5B_1flp, lcg5B_lgcwC, lcg5B_1h1b, lcg5B_lithA, lcg5B_2fal, lcg5B_2hbg, lcg5B_3sdhA, lchkA_1192, lchmA_1ihoA, lcid_1kb5B, lcid_1qa9A, lqipA_1ctqA, lqipA_1hurA, lcjwA_1i21A, lcjxA_1qipA, lckeA_1dekA, lckuA_1eytA, lckv_1gl0A, lcl2A_1qgnH, lclc_1tf4A, lclqA_1nozB
23	lcmbA_1mjoB, lcmoA_1hjbC, lcnt3_1bgc, lcnt3_1cnt2, lcnt3_1lki, lcnv_1llo, lcnzA_1iso, lcof_1cnuA, lcozA_1ihoA, lcozA_2ts1, lcp2A_1eg7A, lcp2A_2nipA, lcpca_1alla, lcpca_2hbg, lcpq_256bA, lcpt_1f4tA, lcqkA_1fc1B, lcqkA_1hxmD, lcqkA_1iakA, lcqqA_1havA, lcqxA_1dlyA, lcr1A_2reb, lcr5A_1e32A, lcr5A_1qcsA, lcr6B_1ehyA
24	lcs6A_1fltX, lcs6A_1tit, lcs8A_1cv8, lcs8A_1qmyA, lcsH_4ctsA, lcsn_1ckiA, lct9A_1gdoA, lctj_1e29A, lctj_2mtaC, lctn_1d2kA, lcto_1bpv, lctqA_1am4D, lctqA_1d4aA, lctqA_1d5cA, lcv8_1dkiA, lcvia_1rpt, lcvjB_2u1a, lcvl_1tca, lcvsc_1ev2E, lcvsc_1fltX, lcvsc_1qfoA, lcvsc_1wit, lcvsc_2ncm, lcwpA_1f15B, lcwvA_1qfhA
25	lcwyA_1bag, lcx8B_1xjo, lcx1A_1cgu, lcx1A_7taa, lcxqA_1b9dA, lcxqA_1c0mC, lcy5A_3ygsP, lcydA_1e6wA, lcydA_1e7wA, lcydA_1enp, lcydA_1eny, lcydA_1fds, lcydA_1fjhB, lcydA_1gegG, lcydA_1h5qL, lcydA_1he2A, lcydA_1oaa, lcydA_1ybvA, lcydA_2ae2B, lcyo_1cxyA, lcyx_1plc, lcyx_2occB, lczfA_1rmg, lcztA_1eut, lcztA_1ulo
26	ld0nA_2vil, ld1dA_1qrjB, ld1gA_1vdrA, ld1gA_3dfr, ld1rA_2if1, ld2fA_1c7nF, ld2fA_1c7oH, ld2iA_1es8A, ld2kA_1ctn, ld2kA_1e9lA, ld2mA_1pjr, ld2nA_1g6oA, ld2sA_1dykA, ld2sA_1sacA, ld2tA_1eoiA, ld2zA_1d2zB, ld2zA_1fadA, ld2zA_1ngr, ld2zB_1fadA, ld2zB_1ngr, ld3bA_1d3bB, ld3bA_1i8fA, ld3bB_1i8fA, ld3gA_2dorA, ld4oA_1j8fA
27	ld4tA_1csyA, ld4tA_1jwoA, ld4vA_1extA, ld4vB_2tnfA, ld5rA_1fpzA, ld5rA_1vhrA, ld5yA_1bl0A, ld6aB_1dm0A, ld6jA_1f48A, ld6jA_1nksA, ld6jA_1nstA, ld6jA_3tmkA, ld7yA_1nhp, ld8jA_1d8kA, ld9cA_1ekuB, ld9eA_1dvjA, ld9eA_1fx6B, ld9eA_1qr7A, ldaaA_1ekfA, ldaaA_1et0A, ldaaA_1i1mC, ldapA_1dssG, ldapA_1hyhA, ldar_1efcA, ldb1A_2lbd
28	ldb1A_3erdA, ldb3A_1bg6, ldb3A_1fds, ldb3A_1he2A, ldb3A_1qrrA, ldbtA_1dvjA, ldbwA_1tmy, ldbwA_3chy, ldceB_1ft1B, ldcfA_1dbwA, ldcfA_3chy, ldciA_1ef8A, ldciA_2dubF, ldfs_1bk0, lddbA_2bidA, lddzA_1i6pA, ldebA_1fe6A, ldekA_1gky, ldekA_3tmkA, ldeoA_1wab, lderA_1a6dA, ldfuP_1feuA, ldg9A_1d2aB, ldg9A_1jfvA, ldgnA_3crd

ตารางที่ 16 กลุ่มข้อมูลที่ 29 - 35

กลุ่มข้อมูล	ชื่อชุดข้อมูล
29	ldhn_1b9IA, ldhpA_1nal3, ldhpA_2tpsA, ldhr_1ybvA, ldhx_1ruxA, ldi0A_1rvv1, ldi1A_1ps1A, ldi6A_1ihcA, ldin_1ei9A, ldjxB_1rlw, ldk4A_1imbA, ldk4A_1inp, ldk5A_1aow, ldkiA_1cv8, ldk1A_1qfxA, ldkuA_1qb7A, ld12A_1hcuA, ld15A_1dusA, ld15A_1fbnA, ld15A_1g6q2, ld15A_1vid, ld15A_1xvaA, ldleA_1qq4A, ldleA_1svpA, ldlyA_1dlwA
30	ldlyA_2gdm, ldlyA_2vhbA, ldmgA_1ffkC, ldmhA_3pchA, ldmhA_3pchM, ldn1A_1dcfA, ldo0A_1e32A, ldo0A_1g8pA, ldokA_1tvxA, ldosA_1zen, ldpe_1jevA, ldpsA_1qghA, ldptA_1ca7A, ldptA_1otgA, ldpuA_1hstA, ldpuA_1qbjA, ldpuA_1smtA, ldqaA_1qaxA, ldqnA_1tc1A, ldqrA_1iatA, ldqrA_2pgi, ldquA_1pymA, ldqwA_1dvjA, ldqyA_1ei9A, ldqyA_1ivyA
31	ldr9A_1i85A, ldrmA_1ew0A, ldrmA_3pyp, ldrw_1gcuA, ldrw_1id1A, ldrw_1ofgA, ldssG_1ceqA, ldssG_1he2A, ldt4A_1vih, ldt6A_1bu7A, ldtyA_2dkb, ldtyA_2gsaA, ldulA_1hq1A, ldun_1dupA, ldusA_1i4wA, ldusA_2dpmA, ldvjA_1fwrA, ldvjA_1ho1A, ldvjA_2tpsA, ldvpA_1elkA, ldwnA_1msc, ldwnA_1qbeA, ldxy_1gdhA, ldxy_1psdA, ldxy_2dl1A
32	ldyKA_1a3k, ldynA_1mai, ldynA_1pls, ldyoA_1ulo, ldz1A_1ap0, ldz3A_1iow, ldz3A_3chy, le0cA_1rhs, le15A_1ctn, le15A_1e9IA, le20A_1g5qA, le2tA_1f13A, le2xA_1qbjA, le32A_1e69B, le3jA_1pedA, le3jA_1qr6A, le4fT_1g99A, le54A_2omf, le54A_3prn, le5dA_5nul, le69B_1f2tA, le69B_1f2tB, le69B_1gajA, le69B_1qhlA, le6bA_1aw9
33	le6bA_1eemA, le6bA_1gnwA, le6wA_1e7wA, le6wA_1enp, le6wA_1fds, le6wA_1fjhB, le6wA_1oaa, le6wA_1ybvA, le70M_1qvbA, le70M_1tr1B, le7wA_1ybvA, le8xA_1e8zA, lea1A_1bu7A, leaf_1c4tA, leaf_3cla, leagA_1smrA, leagA_2er7E, leagA_2rmpA, leagA_3pep, leaiC_1ate, lebmA_1mun, lebuA_1gcuA, lecfB_1gdoA, lecpA_1a2zA, lecpA_1k3fA
34	lecxA_1elqA, lecxA_2gsaA, leduA_1dvpA, leduA_1hg5A, leemA_1gnwA, leemA_1hqoA, leerA_1f6fA, leerB_1iarB, lef1A_1a5r, lef1A_1gg3C, lefdN_1qguA, lefuB_1tfe, lefvA_1efpA, lefvA_1mjhA, leg7A_1hyqA, lehaA_1uok, lehyA_1qo7A, lehyA_2bce, lei5A_3pte, lei9A_1cr6B, lei9A_1cvl, leia_1d1dA, leiwA_1egaA, lej0A_1dhr, lej0A_1fbnA
35	lejeA_1axj, lejeA_1i0rA, lejha_1ap8, lekeA_2rn2, lekfA_1et0A, lekrA_1mla, lelqA_2dkb, lema_1ggxA, lemsA_1erzA, lenp_1ej0A, lenp_1eny, lenp_1qsgA, lenp_1ybvA, leny_1ybvA, leokA_1d2kA, lep2B_1fnc, lepaA_1bj7, lepaA_1mup, lepaA_1np1A, leq2A_1xel, leq3A_1pinA, leqkA_1cewI, leqkA_1molA, lerv_1thx, lerzA_1f89A

ตารางที่ 17 กลุ่มข้อมูลที่ 36 - 42

กลุ่มข้อมูล	ชื่อชุดข้อมูล
36	lesc_1eny, lesl_1h8uA, lesl_1qddA, lesl_1qo3C, lesl_1rtm1, letb1_1gkeC, leteA_1hmcB, letpA_1fcdC, leu8A_4mbp, leuhA_1eyyA, leumA_1qghA, leumA_2fha, leut_2sli, leut_3sil, levhA_1i7aC, levhA_1qc6A, levsA_1lki, lew4A_1dlxA, lewaA_1flp, lewaA_1lithA, lewaA_2fal, lewxA_1gp1A, lewxA_1jfuA, lex1A_3chy, leyg_1xbd
37	lextA_1d4vA, lexA_1hmcB, leyrA_1fwyA, leyrA_1ga8A, leyrA_1h7eA, leyyA_1ad3A, lezvB_1be3B, lf13A_1g0dA, lf2dA_1oasA, lf2nA_1smvA, lf2nA_2tbvA, lf37A_1b9yC, lf3yA_1g0sA, lf4tA_1egyA, lf4tA_1phd, lf4tA_1rom, lf53A_1g6eA, lf5qB_1guxB, lf5wA_1neu, lf5xA_1foeA, lf7cA_1pbwA, lf7cA_1tx4A, lf94A_1cds, lf94A_3ebx, lf97A_1qfoA
38	lf9aA_1hybA, lfaoA_1dynA, lfaoA_1ef1A, lfaoA_1fhoA, lfaoA_1fhxA, lfaoA_1mai, lfaoA_1rrpB, lfc6A_1pdr, lfc6C_1e29A, lfchA_1hxiA, lfdo_1dmr, lfdr_1a8p, lfdr_1fnc, lfdr_1qfjA, lfds_1he2A, lfepA_1by5A, lff9A_1gcuA, lff9A_1id1A, lff9A_2scuA, lffkE_1e7kB, lffkC_1ocrC, lfggA_1fgxA, lfggA_1fr9A, lfggA_1ga8A, lfgkA_1csn
39	lfgkA_1fgiB, lfgkA_1ia8A, lfgxA_1foaA, lfgxA_1g8oA, lfgxA_1ga8A, lfh6A_1fh6G, lfhoA_1pls, lfhoA_1qqgA, lfhuA_1mucA, lfi2A_1caxB, lfi2A_1dzaA, lfit_1kpf, lfj7A_1cvjB, lfjcA_1cvjB, lfjfl_1fjfQ, lfjhB_1eq2A, lfjjA_1qouA, lfknA_1hvc, lfknA_1lywA, lfknA_1smrA, lfknA_2rmpA, lfloA_1ae9A, lfloA_1aihA, lfip_1h97A, lfip_1h1b
40	lfip_1lithA, lfip_2fal, lfip_3sdhA, lfltX_1qa9A, lfltX_1tit, lfltX_2ncm, lfm5A_1htn, lfm5A_1qddA, lfm5A_1qo3C, lfmb_1b11A, lfmb_2hpeA, lfmb_2rmpA, lfmk_1hcl, lfmtA_2gar, lfnc_1fb3B, lfnc_1fdr, lfnc_1qfjA, lfnc_2pia, lfo5A_1hyuA, lfofA_1e25A, lfppA_1dusA, lfppA_1fp1D, lfpaA_1pty, lfpaA_1vhrA, lfqtA_1rfs
41	lfqtA_1rie, lfqvD_1vcbB, lfr9A_1fgxA, lfr9A_1h7eA, lfr9A_1i52A, lfrb_1qrqA, lfrpA_1dk4A, lfshA_1hstA, lfi9A_1qbjA, lfi9A_2cgpA, lfukA_1qvaA, lfvaA_1qd1A, lfvaA_1bed, lfvaA_1bv1, lfwkA_1fi4A, lfxd_1vjw, lfxkA_1fxkC, lfxkC_1fxkA, lfzgD_2eboA, lg0rA_1qg8A, lg10A_1ckv, lg1eB_1e91A, lg24A_1qs1A, lg4uS_1ytw, lg4uS_2shpA
42	lg55A_6mhtA, lg5rA_1g6oA, lg5rA_1g8yA, lg5rA_2reb, lg61A_1g62A, lg6gA_1qu5A, lg6q2_1f31A, lg6q2_1vid, lg6q2_2admA, lg7eA_1prxA, lg8fA_1coza, lg8fA_1f9aA, lg8jA_2napA, lg8jB_1fqtA, lg8jB_1rfs, lg8jB_1rie, lg8oA_1qg8A, lg8qA_1g8qB, lg8yA_1g5rA, lga3A_1rcb, lga8A_1fgxA, lga8A_1g8oA, lga8A_1qg8A, lgakA_1lis, lgaxA_1ile

ตารางที่ 18 กลุ่มข้อมูลที่ 43 - 49

กลุ่มข้อมูล	ชื่อชุดข้อมูล
43	lgc1G_1g9nG, lgca_1qpzA, lgca_1tlfA, lgca_2dri, lgceA_3pte, lgci_1sud, lgdhA_1psdA, lgefA_1hh1A, lgen_1hxn, lgen_1pex, lggqA_1flmC, lggxA_1h4uA, lghj_1iyu, lghqB_1e5gA, lgkxA_1id0A, lgky_1nksA, lgky_1qhsA, lgky_3tmkA, lgnwA_1ljrA, lgnwA_1pgtA, lgp1A_1erv, lgp1A_1jfuA, lgp1_1bu8A, lgpma_1qdlB, lgr2A_1ii5A
44	lgsa_2hgsA, lgtxA_2dkb, lgtxA_2gsaA, 1h70A_1jdw, 1h75A_1kte, 1h7wA_2dorA, 1h8cA_1c1yB, 1h8cA_1ubi, 1h8uA_1htn, 1h8uA_1qddA, 1h8uA_1rtm1, 1h97A_1ithA, 1h97A_2fal, 1h97A_2hbg, 1h97A_3sdhA, 1h9jA_1g291, 1ha1_1f9fA, 1ha1_1hd0A, 1ha1_1jmtA, 1ha1_2u2fA, 1han_1cxA, 1han_1eirA, 1han_1jc4A, 1havA_1cqqA, 1hcl_1csn
45	1hd2A_1qq2A, 1hdmB_1b3jA, 1he2A_1hyhA, 1he2A_1qrrA, 1hg3A_1nsj, 1hg3A_2tpsA, 1hgeB_1flcB, 1hjp_1bvsa, 1h1b_1dlyA, 1h1b_2fal, 1hmcB_1jli, 1hmt_1lfo, 1hn9A_1bi5A, 1hnoA_1nzyA, 1ho1A_1nsj, 1hp4A_1qba, 1hqoA_1g6wA, 1hqoA_1gnwA, 1hqoA_1ljrA, 1hqoA_1pgtA, 1hs6A_1hyt, 1hssA_1bea, 1hssA_1hyp, 1htn_1ixxB, 1htn_1rdo1
46	1htn_1rtm1, 1huuA_1hueA, 1huw_1a22A, 1hv8A_1c9kB, 1hwyA_1b26F, 1hwyA_1k89, 1hx1B_2a3dA, 1hxmA_1tvdA, 1hxmD_1tvdA, 1hyhA_1hygB, 1hyhA_1ldnF, 1hyhA_2cmd, 1hyhA_3ldh, 1hyt_1ezm, 1hyt_1hs6A, 1hyuA_1tde, 1i1jA_1ckaA, 1i21A_1i12D, 1i21A_1nmtA, 1i21A_1yghA, 1i4wA_1vid, 1i4wA_2dpmA, 1i5nA_1qspA, 1i69A_1pda, 1i6wA_1jfrA
47	1i6wA_1tca, 1iae_1kuh, 1iakA_1hdmA, 1iakA_1tlk, 1iarB_1cto, 1ibjA_1bs0A, 1ibjA_1cs1A, 1ibzA_1kcw, 1ibzA_1plc, 1ibzA_1rcy, 1iciA_1nbaA, 1id1A_1eq2A, 1id1A_1ofgA, 1ifa_1huw, 1ifa_1rmi, 1ig0A_1ig3A, 1igtB_1fe8J, 1ihmA_1a6cA, 1ihp_1bif, 1ihp_1qfxA, 1ii5A_1lst, 1ijqA_1qlgA, 1im3D_1tvdA, 1imbA_1jp4A, 1imbA_1qgxA
48	1iow_1ehiB, 1iow_1gpmA, 1iq3A_2cblA, 1irp_2irtA, 1iso_1cnzA, 1itbB_1kb5B, 1itbB_1tit, 1itbB_1tlk, 1itbB_1wit, 1itbB_2ncm, 1ithA_1cg5B, 1ithA_2fal, 1ithA_3sdhA, 1ixh_1pot, 1ixxB_1eggA, 1iyu_1dczA, 1iyu_1ghj, 1j9yA_1oyc, 1jdc_1amy, 1jfiB_1tafA, 1jfiB_1tafB, 1jfuA_1prxA, 1jhjA_1xnaA, 1jhnA_1led, 1jk0A_1xikA
49	1jkmA_1evqA, 1jmtA_2u2fA, 1jn5A_1ounA, 1jn5B_1opy, 1jotA_1c3nA, 1jotA_1vmoA, 1jp4A_1qgxA, 1jxpA_1svpA, 1k89_1bw9A, 1kb5B_1fo0B, 1kb5B_1nkr, 1kb5B_1qa9A, 1kb5B_1qfoA, 1kb5B_1tlk, 1kb5B_1vcaA, 1kb5B_2ncm, 1kit_2sli, 1kjs_1c5a, 1krs_1b8aA, 1krs_1ly1A, 1ksiB_1spuB, 1kuh_1hfc, 1kuh_1iae, 1kum_1cx1A, 1kwaA_1kwaB

ตารางที่ 19 กลุ่มข้อมูลที่ 50 - 56

กลุ่มข้อมูล	ชื่อชุดข้อมูล
50	lkwaA_1pdr, lkwaA_1qauA, llam_1ecpA, llarA_1d5rA, llarA_2shpA, llay_1at3A, llbd_2lbd, llckA_1griA, llcl_1sacA, llea_1qbjA, llfb_2hddA, llfdA_2rgf, llis_1gakA, lljrA_1pd21, llki_1exzA, llla_1hc2, llpbA_1pcn, llst_1nnt, llt3A_1bcpA, llt3A_1xtcA, llvk_1br2F, llvl_1geuB, llxa_1qreA, llxa_2xat, lmai_1pls
51	lmaz_1fl6A, lmaz_1g5jA, lmaz_2bidA, lmfa_1neu, lmfmA_1xsoB, lmfmA_2apsB, lmgtA_1sfe, lmh1_2cmd, lmmA_1c7uB, lmoq_1bvyF, lmpgA_1mun, lmpgA_2abk, lmrj_1qcjB, lmroA_1mroB, lmroB_1e6vE, lmrp_1d9yA, lmsc_1qbeA, lmspA_1qpxA, lmtb_1mhyB, lmtb_1qq8A, lmucA_2mnr, lmugA_1bjt, lmugA_3eugA, lmun_1ebmA, lmup_1bj7
52	lmup_1qftA, lnar_1d2kA, lnbaA_1yacA, lnbcA_1g43A, lnbcA_1tf4A, lndoA_1rie, lnedA_1g3iJ, lnedA_1pmaP, lneu_1kacB, lneu_1tvdA, lnflA_1wer, lngl_1j8yF, lnksA_1nstA, lnksA_1qhiA, lnksA_1shkA, lnmtA_1qsmA, lnnt_1dsn, lnox_1bkjA, lnkp_1nueD, lnseA_1qomB, lnsb_256bA, lnukA_1ulo, lnwpA_1qhqA, lnzyA_1dubE, loaa_1ybvA
53	loasA_1tdj, lobpA_1mup, locrC_1qleC, lopc_1qqiA, lopy_3stdA, lorc_4croE, lotcB_1quqA, lovaA_1sek, lovaA_2achA, loyc_2tmdA, lp35A_1i4eA, lpauB_1qduL, lplc_1czfA, lplc_1qcxA, lpd21_1pgtA, lpdgA_1vpfA, lpd_1kwaA, lpd_2pdzA, lpedA_1bxzD, lpedA_2ohxA, lpgtA_1gsdB, lpgtA_1gumH, lpii_1a53, lpjr_1hv8A, lpjr_1qvaA
54	lpjr_1uaaA, lplc_1rcy, lplc_2occb, lplq_1axcE, lplq_1ge8A, lplq_2polA, lpls_1foeA, lpls_1rrpB, lpmaA_1nedA, lpmaA_1pmaP, lpne_3nul, lpot_1lst, lpoxA_1bfd, lpreA_1bcpB, lprs_1g6eA, lprtF_2bosA, lprtF_3chbD, lpscA_1b5tA, lpsrA_1c1l, lpsrA_1qlsA, lpsrA_4icb, lptf_2hid, lpty_1ytw, lpty_2shpA, lpvc4_1ar94
55	lpvl_7ahlA, lqa9A_1qfoA, lqa9A_1tcrA, lqa9A_1tit, lqa9A_1tlk, lqa9A_1wit, lqa9A_2ncm, lqazA_1cem, lqbhA_1g73D, lqbjA_2cgpA, lqbzC_2ezoA, lqdlB_1i7qB, lqdlB_1tmy, lqfjA_1fnc, lqfoA_1tit, lqfoA_1tvdA, lqfoA_1vcaA, lqfoA_1wit, lqfoA_2fcbA, lqfoA_2ncm, lqghA_1bcfA, lqghA_2fha, lqguB_3csuA, lqgvA_1b9yC, lqgvA_1bjx
56	lqgwA_1qgwB, lqgwC_1allA, lqgwC_1cpcL, lqhaA_1g99A, lqhiA_1qhsA, lqhqa_1bqk, lqhsA_3tmkA, lqhuA_1ck7A, lqhvA_1nobE, lqipA_1cxA, lqipA_1fa5B, lqj2A_1ffvA, lqj2B_1ffvE, lqj2C_1ffuF, lqj4A_1ei9A, lqjA_1qo8D, lqk3A_1gca, lqk3A_1nulA, lqksA_1n50A, lqlaA_1fumA, lqlaB_1clyB, lqndA_1c44A, lqniA_1libzA, lqnoA_1eceA, lqnxA_1cfe

ตารางที่ 20 กลุ่มข้อมูลที่ 57 - 63

กลุ่มข้อมูล	ชื่อชุดข้อมูล
57	lqo0D_1ybvA, lqo2A_1thfD, lqo3C_1rtm1, lqorA_2ohxA, lqoxN_1e70M, lqpxA_1bf8, lqpzA_1tlfA, lqpzA_2dri, lqq4A_1arb, lqq4A_1svpA, lqq4A_2hrvA, lqq8A_1j77A, lqqp1_2mev1, lqqp2_1pov1, lqqp3_1mec3, lqqp3_1pov1, lqqp3_2mev1, lqr4B_1f97A, lqrb_1ak4C, lqs1A_1g24A, lqsmA_1yghA, lqtrA_1hlgA, lqtrA_1qfmA, lqu5A_1ygs, lqu9A_1qd9C
58	lqpA_1yaiA, lquqA_1quqB, lr2fA_1xikA, lr2fA_1xsm, lr69_1b0nA, lrcb_1lki, lrcb_3inkC, lrcf_1d04A, lrcf_5nul, lrcy_1qhqA, lrcy_2cuaA, lrec_1bjfB, lrfs_1rie, lrhoA_1ft3A, lrhs_1e0cA, lrie_1ezvE, lrom_1jipA, lrpA_2dorA, lrsy_1a25B, lrthA_2m2, lrtm1_1qo3C, lrzl_1fk3A, lsacA_1b09C, lsat_1hfc, lsayA_1hzzB
59	lsbp_1wod, lscjB_1b64, lseiA_1qd7G, lsek_1dvmD, lsek_1psi, lsfe_1eh7A, lshcA_1x11A, lshkA_1ng1, lshkA_3tmkA, lshsA_1ejfA, lskyB_1a5t, lskyE_1e32A, lskyE_1skyB, lsluA_1fi8E, lsmrA_2rmpA, lsmvA_1c8nC, lsmvA_2tbvA, lsmvA_4sbvC, lsro_1mjc, lstfl_1cewI, lsur_2nsyA, lsvpA_5ptp, lsvy_1svq, lsvy_2vil, lswuA_2aviA
60	lt1dA_1buoA, ltaq_5ktqA, ltaxA_1xas, ltbge_1tbgF, ltc1A_1qk3A, ltcrA_1ahl, ltcrA_1qsfD, ltfb_1c9bA, ltf_1efuB, ltgxA_1cdtA, ltheB_1ef7A, ltheB_1qmyA, ltig_2ifeA, ltit_1wit, ltit_2ncm, ltlfA_2dri, ltlk_1wit, ltlk_1wwcA, ltrb_1f6mF, ltul_1dun, ltvdA_1bzqN, ltvdA_1ivlB, ltvdA_1qfwL, ltvxA_1a15A, ltx4A_1f7cA
61	ltyfA_1nzyA, lu2fA_1fjcA, lu9aA_1c4zD, lu9aA_1i7kB, lubpC_1a4mA, luch_1cmxA, lulo_1cx1A, luok_1bvzA, lurnA_1u2fA, lurnA_2u1a, luteA_1qhwA, luteA_4kbpA, lvcaA_2fcbA, lvdrA_3dfr, lvfrA_1licuB, lvfrA_1nox, lvfyA_1hyiA, lvid_2admA, lvid_2dpmA, lvid_3mag, lvpfA_1pdgA, lvpsA_1dzlA, lvsgA_2vsgA, lwab_1deoA, lwdcC_1br4B
62	lwer_1nf1A, lwgjA_2prd, lwit_2ncm, lwwcA_1he7A, lx11A_2nmbA, lxbd_1ayoA, lxel_1i3nA, lxikA_1qghA, lxikA_1xsm, lxvaA_2dpmA, lxvaA_2ercA, lxwl_1t7pA, lxwl_2kfzA, lxxaA_1b4aA, lybvA_1g6q2, lybvA_2ae2B, lyer_1b63A, lyge_1byt, lyghA_1bob, lyghA_1i21A, lyghA_1qsrA, lygs_1devA, lytbA_1qnaB, lzin_5ukdA, lzpdA_1qpbB
63	lzxq_1ic1A, 2a3dA_1fpoA, 2abk_1ebmA, 2admA_3mag, 2afpA_1ixxB, 2bbkH_1qfmA, 2bbvA_1f2nA, 2bce_1acl, 2bce_1qonA, 2bopA_1ris, 2btvA_2btvB, 2cba_1koqB, 2cblA_1psrA, 2cbp_1f56C, 2ccyA_256bA, 2cmd_1mldD, 2cmd_3ldh, 2cpl_1c5fC, 2cpl_2nul, 2dhqA_1c9kB, 2dkb_2gsaA, 2dkb_2oatC, 2dorA_1ep3A, 2dri_8abp, 2eboA_1mof

ตารางที่ 21 กลุ่มข้อมูลที่ 64 - 68

กลุ่มข้อมูล	ชื่อชุดข้อมูล
64	2er7E_2rmpA, 2ercA_3mag, 2fal_1ebt, 2fal_3sdhA, 2fcbA_1iisC, 2fcbA_2ncm, 2fnbA_1bpv, 2gar_1bxkA, 2gdm_1a6m, 2gdm_2vhbA, 2gmfA_1f6fA, 2hbg_2vhbA, 2hbg_3sdhA, 2hddA_1b8iA, 2hdhA_3hdhB, 2hrvA_1agjA, 2hrvA_5ptp, 2hts_2irfG, 2i1b_1iltA, 2i1b_1iraX, 2i1b_2ila, 2ifl_1d1rA, 2ila_1abrB, 2ila_1hce, 2ilk_1vlk
65	2ibd_2prgA, 2masA_1ezrA, 2mbr_1qltA, 2mevl_1tmf1, 2nlrA_1h8vF, 2nmbA_1shcA, 2nsyA_1sur, 2occB_1qleB, 2ohxA_3hudA, 2omf_1osmC, 2omf_2por, 2pgd_1pgjA, 2pii_1cc8A, 2pkaA_1a7s, 2pkaA_1aksA, 2por_3prm, 2prgA_3erdA, 2pth_1c8bA, 2qwc_1b9vA, 2qwc_1nsdB, 2reb_1cr1A, 2reb_1g19A, 2rmpA_1avfA, 2rn2_1ekeA, 2scuA_1drw
66	2scuA_1eudA, 2scuB_1eucB, 2shpA_1i9sA, 2tbvA_1pov1, 2tct_1a6i, 2tgi_1es7C, 2thiA_1eu8A, 2tnfA_1d4vB, 2tpsA_1qpoA, 2tysA_1pii, 2tysB_1tdj, 2u2fA_1u2fA, 2xat_3tdt, 3bct_1g3jC, 3chbD_1prtF, 3chy_1b00B, 3chy_1tmy, 3chy_4tmyB, 3cla_1nocB, 3crd_3ygsP, 3erdA_1dkfA, 3erdA_1ereB, 3eugA_1lauE, 3inkC_1irl, 3ldh_1hlpA
67	3lzt_1gbzA, 3nul_1acf, 3nul_1fil, 3nul_1ifqA, 3pchM_3pchA, 3pte_1pmd, 3pte_1skf, 3pyp_1drmA, 3sdhA_1sctF, 3seb_1enfA, 3stdA_2std, 3tgl_1dt5A, 3ulla_1kawB, 3ulla_1prtF, 3vub_2vubH, 4aahA_1flgB, 4bcl_1ksaA, 4crxA_1floA, 4icb_1mho, 4mbp_1eljA, 4mbp_1eu8A, 4mbp_1ezpA, 4mbp_1gggB, 4pgaA_1ho3A, 4uagA_1eehA
68	5hpgA_2pk4, 5nul_1akt, 5nul_1j9gA, 5tmpA_3tmkG, 6prcL_6prcM, 7taa_2taaA, 8fabA_1tetH

รวมทั้งหมด 68 กลุ่มข้อมูลมี 1682 ชุดข้อมูลที่ใช้ในการทดลอง

เนื่องจากการหาคูฐานข้อมูล BALiBASE และ PREFAB ในแต่ละชุดที่เทียบเรียงมีจำนวนและความยาวของลำดับข้อมูลไม่มากซึ่งทำให้ผลความแตกต่างทางด้านเวลาที่ใช้ในการประมวลผลในแต่ละชุดข้อมูลไม่เด่นชัด ดังนั้นจึงได้มีการวัดผลทดลองการเทียบเรียงกลุ่มลำดับข้อมูลโปรตีนจากฐานข้อมูล NCBI ที่มีจำนวนข้อมูลตั้งแต่ 10 – 500 ลำดับ และมีความยาวของลำดับเฉลี่ยประมาณ 280 ตัวอักษร เพื่อศึกษาแนวโน้มและเปรียบเทียบเวลาที่ใช้ในการเทียบเรียงของทั้ง 2 โปรแกรมเมื่อจำนวนของลำดับที่เทียบเรียงเพิ่มขึ้นอีกด้วย

3. การวัดความถูกต้องของการเทียบเรียง

ความถูกต้องของผลการเทียบเรียงในงานวิจัยนี้วัดผลโดยพิจารณาจากค่าคะแนนความถูกต้อง ซึ่งสำหรับฐานข้อมูล BALiBASE วัดผลจากค่าคะแนน SP และค่าคะแนน TC และสำหรับฐานข้อมูล PREFAB วัดผลจากค่าคะแนน Q โดยค่าคะแนนแต่ละแบบมีความหมายดังนี้

3.1 ค่าคะแนน SP (Sum of Pair score) คำนวณโดยโปรแกรมของคณะวิจัยผู้สร้างฐานข้อมูล BALiBASE เป็นค่าคะแนนที่ได้จากการหาอัตราส่วนระหว่างจำนวนอักขระของส่วนจำเพาะที่มีตำแหน่งถูกต้องของผลการเทียบเรียงของโปรแกรมเปรียบเทียบกับจำนวนอักขระส่วนจำเพาะทั้งหมดของผลการเทียบเรียงอ้างอิง ค่าคะแนน SP ใช้วัดความถูกต้องของผลการเทียบเรียงฐานข้อมูล BALiBASE บ่งบอกถึงความถูกต้องของผลการเทียบเรียงกลุ่มลำดับข้อมูลเมื่อเปรียบเทียบกับผลเทียบเรียงอ้างอิงซึ่งเทียบเรียงโดยผู้เชี่ยวชาญ ซึ่งค่าคะแนน SP จะมีค่าระหว่าง 0.0-1.0 โดยที่ค่าคะแนนจะเพิ่มขึ้นตามความถูกต้องของผลการเทียบเรียง

3.2 ค่าคะแนน TC (Total Column score) คำนวณโดยโปรแกรมของคณะวิจัยผู้สร้างฐานข้อมูล BALiBASE เป็นค่าคะแนนที่ได้จากการหาอัตราส่วนระหว่างจำนวนของสดมภ์ในส่วนจำเพาะที่ถูกต้องของผลการเทียบเรียงของโปรแกรมเปรียบเทียบกับจำนวนสดมภ์ในส่วนจำเพาะทั้งหมดของผลการเทียบเรียงอ้างอิง ค่าคะแนน TC ใช้วัดความถูกต้องของผลการเทียบเรียงฐานข้อมูล BALiBASE บ่งบอกถึงความถูกต้องของผลการเทียบเรียงในแนวสดมภ์เมื่อเปรียบเทียบกับผลเทียบเรียงอ้างอิงซึ่งเทียบเรียงโดยผู้เชี่ยวชาญเช่นกัน ค่าคะแนน TC จะมีค่าระหว่าง 0.0-1.0 โดยที่ค่าคะแนนจะเพิ่มขึ้นตามความถูกต้องของผลการเทียบเรียงเช่นเดียวกับค่าคะแนน SP

3.3 ค่าคะแนน Q คำนวณโดยโปรแกรมของคณะวิจัยผู้สร้างฐานข้อมูล PREFAB เวอร์ชัน

4.0 ไว้สำหรับวัดผลการเทียบเรียงฐานข้อมูล PREFAB ค่าคะแนน Q คำนวณโดยเริ่มจากนำผลการเทียบเรียงของโปรแกรมไปหาตำแหน่งที่ตรงกันทีละคู่ลำดับ และทำเช่นเดียวกันนี้กับผลการเทียบเรียงอ้างอิง หลังจากนั้นจึงเอาข้อมูลตำแหน่งตรงกันที่หาได้ของแต่ละคู่จากผลการเทียบเรียงของโปรแกรม และผลการเทียบเรียงอ้างอิงมาเทียบกัน โดยคำนวณเป็นค่าเฉลี่ยของอัตราส่วนระหว่างจำนวนของตำแหน่งที่ถูกต้องของผลการเทียบเรียงของโปรแกรม กับจำนวนตำแหน่งที่ตรงกันของแต่ละคู่ของการเทียบเรียงอ้างอิง โดยค่าคะแนน Q จะมีค่าระหว่าง 0 ถึง 1 เช่นเดียวกับค่าคะแนน SP และ TC ของฐานข้อมูล BALiBASE

โปรแกรมสำหรับการคำนวณค่าคะแนน SP และค่าคะแนน TC ชื่อโปรแกรม Bali_score.c สามารถหาได้จาก http://bips.u-strasbg.fr/fr/Products/Databases/BAlIbASE/bali_score.c และโปรแกรมสำหรับการคำนวณค่าคะแนน Q สามารถหาได้จาก http://www.drive5.com/qscore/qscore_src.tar.gz

ผลและวิจารณ์

ระบบที่ใช้ในงานวิจัยนี้ เราใช้เครื่องคอมพิวเตอร์ Athlon XP 2000+ ที่มี CPU Clock rate 1.67 GHz หน่วยความจำหลัก 512 MB ในส่วนของการทดลอง เราใช้โปรแกรม ClustalW และ โปรแกรม SSFA เทียบเรียงชุดลำดับข้อมูลของฐานข้อมูล BALiBASE และ PREFAB เพื่อเปรียบเทียบประสิทธิภาพทางด้านเวลาที่ใช้ประมวลผล และค่าความถูกต้องของผลลัพธ์ที่เทียบเรียงได้ โดยที่ ClustalW ก็คือผลที่ได้จากโปรแกรม ClustalW และ SSFA ก็คือผลที่ได้จากโปรแกรม SSFA

1. ผลการเทียบเรียงชุดข้อมูลของฐานข้อมูล BALiBASE

การเทียบเรียงชุดข้อมูลของฐานข้อมูล BALiBASE วัดผลประสิทธิภาพทางด้านเวลาด้วยเวลารวมในการเทียบเรียงทุกชุดของฐานข้อมูล เนื่องจากแต่ละชุดข้อมูลของฐานข้อมูล BALiBASE มีจำนวนข้อมูลไม่มากส่งผลให้เวลาในการประมวลผลของแต่ละชุดใช้เวลาน้อย และวัดประสิทธิภาพได้ยาก ส่วนในด้านความถูกต้องพิจารณาจากค่าคะแนน SP และค่าคะแนน TC ทั้งของแต่ละชุดข้อมูล และค่าคะแนนเฉลี่ยของแต่ละชุดอ้างอิง

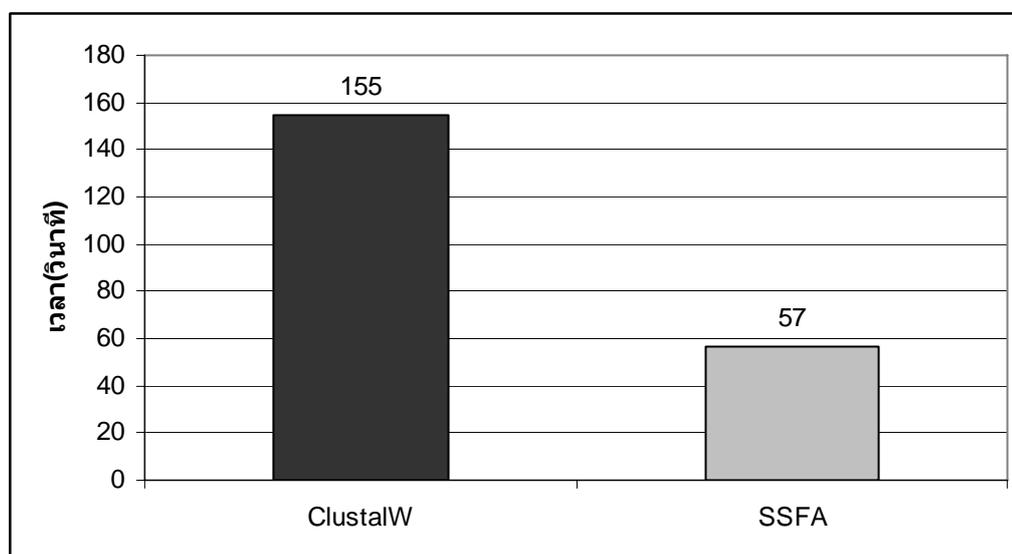
ภาพที่ 17 แสดงเวลาที่ใช้ในการเทียบเรียงลำดับข้อมูลชุดอ้างอิงทั้ง 5 ชุด แสดงให้เห็นว่าที่โปรแกรม ClustalW ใช้เวลารวมในการประมวลผลมากกว่าโปรแกรม SSFA อย่างเห็นได้ชัด เนื่องจากเวลาที่ใช้ในขั้นตอนที่ 1 ของการเทียบเรียงลดลง ซึ่งเวลาที่โปรแกรม ClustalW ใช้ นั้นมากกว่าโปรแกรม SSFA ถึง 2.72 เท่า

ภาพที่ 18 - ภาพที่ 31 แสดงการเปรียบเทียบค่าคะแนน SP และค่าคะแนน TC ของผลลัพธ์ที่ได้จากการเทียบเรียงด้วยโปรแกรม ClustalW และ โปรแกรม SSFA สำหรับแต่ละชุดข้อมูล และภาพที่ 32 และภาพที่ 33 แสดงการเปรียบเทียบค่าคะแนน SP และค่าคะแนน TC เฉลี่ยของแต่ละชุดอ้างอิงเมื่อเทียบเรียงด้วยโปรแกรม ClustalW และ โปรแกรม SSFA

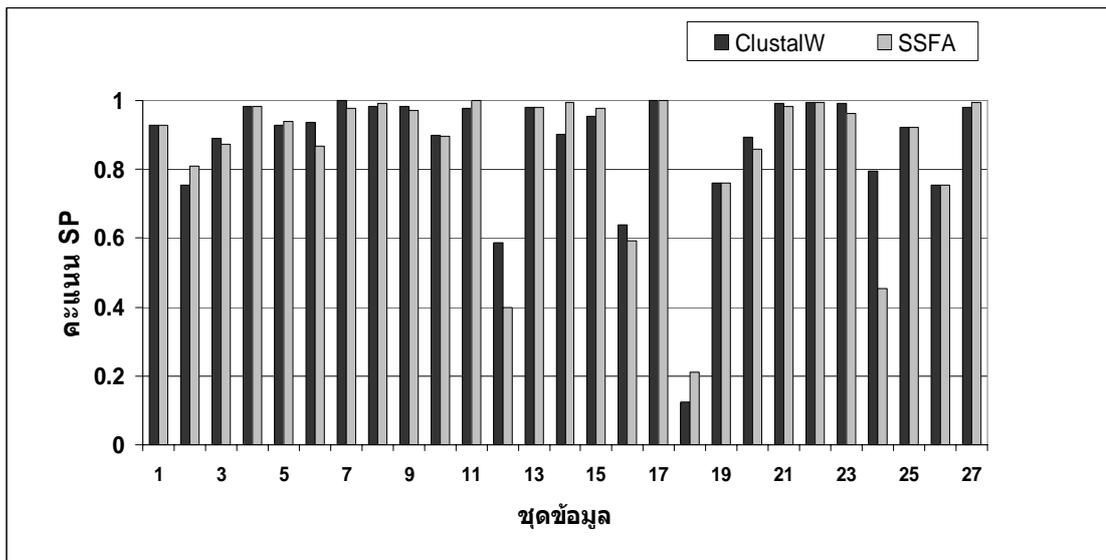
เมื่อพิจารณาค่าคะแนน SP ในภาพคู่ (ภาพที่ 18 – ภาพที่ 32) แสดงให้เห็นว่าโปรแกรม SSFA และ โปรแกรม ClustalW ให้ค่าคะแนน SP มีแนวโน้มไปในทิศทางเดียวกัน และมีค่าคะแนนใกล้เคียงกันอยู่ในระดับที่ยอมรับได้ โดยเฉพาะชุดอ้างอิง: 5 ซึ่งเป็นชุดข้อมูลที่มีการแทรก

ภายในโปรแกรม SSFA ได้ค่าคะแนน SP ที่ดีกว่าโปรแกรม ClustalW ในหลายๆชุดข้อมูล และโปรแกรม SSFA มีค่าคะแนน SP เฉลี่ยของชุดอ้างอิง: 5 มากกว่าโปรแกรม ClustalW เล็กน้อย

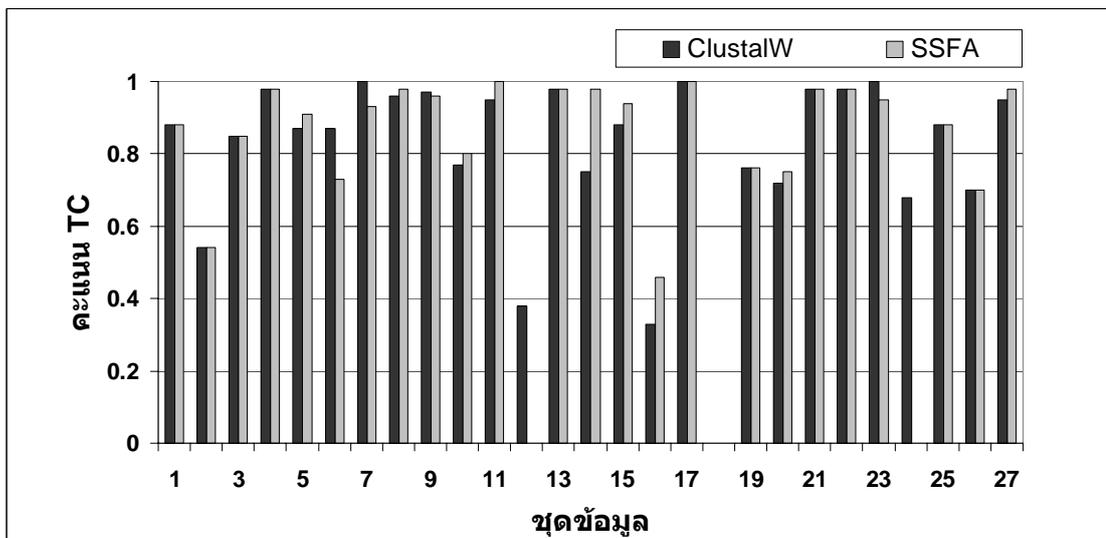
เมื่อพิจารณาค่าคะแนน TC ซึ่งเป็นค่าคะแนนความถูกต้องของสคมภ์เมื่อเปรียบเทียบกับผลลัพธ์อ้างอิง ในภาพคู่ (ภาพที่ 19 – ภาพที่ 33) แสดงให้เห็นว่าโปรแกรม SSFA และโปรแกรม ClustalW ให้ค่าคะแนน TC ที่ใกล้เคียงกันอยู่ในระดับที่ยอมรับได้เช่นเดียวกัน ยกเว้นชุดอ้างอิง: 4 ซึ่งโปรแกรม SSFA ให้ค่าคะแนน TC น้อยกว่าโปรแกรม ClustalW มากพอสมควร เกิดจากการคำนวณค่าคะแนน TC สนใจเฉพาะความถูกต้องในส่วนจำเพาะ และชุดข้อมูลที่ 4 มีลักษณะการมีข้อมูลขนาดใหญ่เพิ่มเติมที่ตอนต้นหรือตอนท้ายของข้อมูล ซึ่งลักษณะข้อมูลแบบชุดข้อมูลที่ 4 นี้ทำให้ผลการเทียบเรียงของโปรแกรม SSFA ได้ส่วนจำเพาะที่กระจายไปในข้อมูลขนาดใหญ่ที่เพิ่มเติมเข้ามา ทำให้ผลการเทียบเรียงแตกต่างจากผลการเทียบเรียงอ้างอิงมากกว่าโปรแกรม ClustalW แต่หากพิจารณาค่าคะแนน SP ของชุดอ้างอิง: 4 ที่มีค่าใกล้เคียงกัน สามารถสรุปได้ว่าโปรแกรม SSFA สามารถเทียบเรียงชุดอ้างอิง: 4 ให้ผลลัพธ์ที่มีความถูกต้องในระดับที่ยอมรับได้ แต่ผลลัพธ์มีความแตกต่างจากผลลัพธ์อ้างอิงมากกว่าโปรแกรม ClustalW



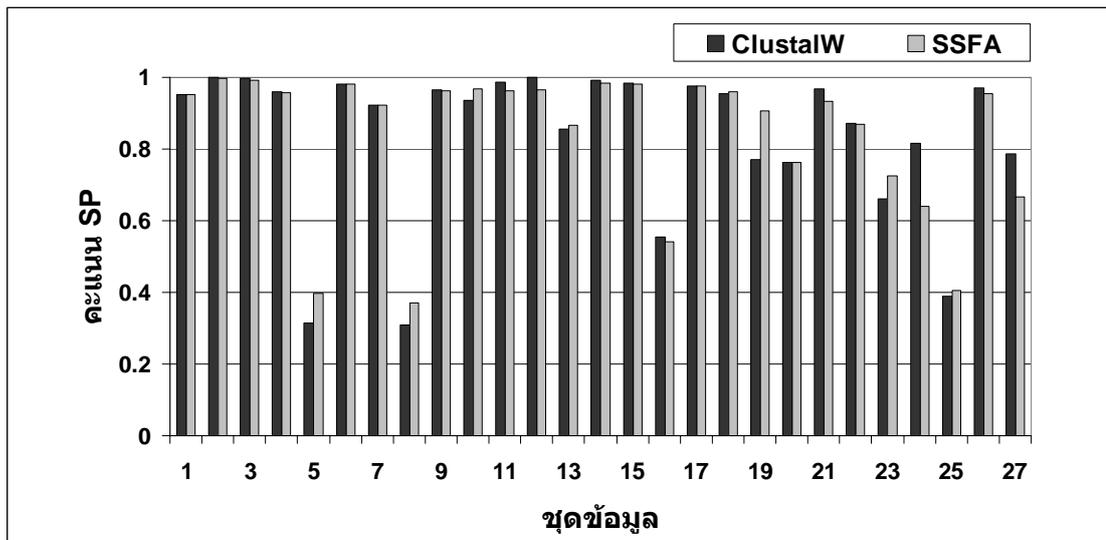
ภาพที่ 17 เวลาที่โปรแกรม ClustalW และโปรแกรม SSFA ใช้ในการเทียบเรียงกลุ่มลำดับข้อมูลของชุดอ้างอิงทั้ง 5 ชุด



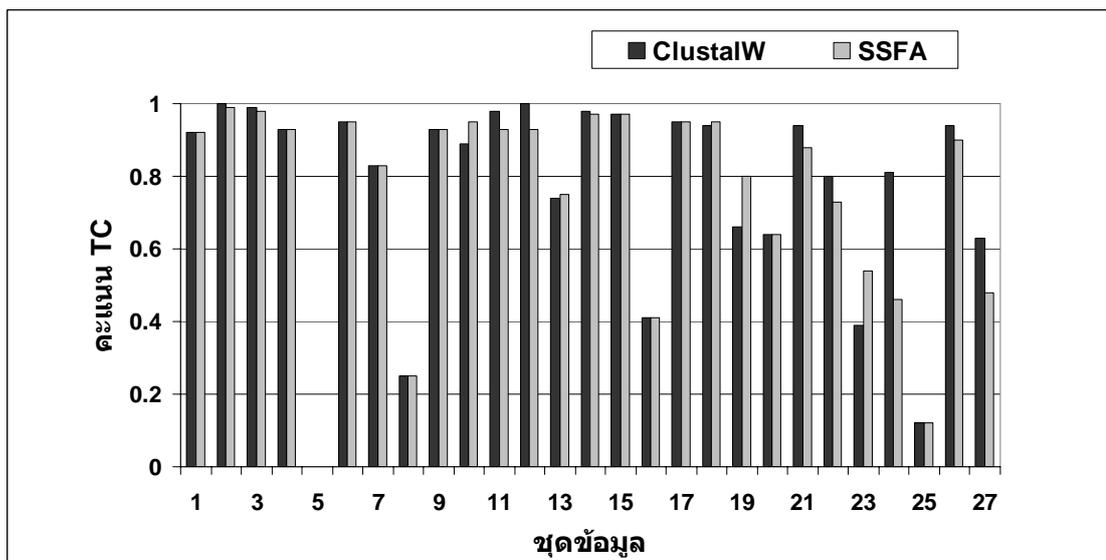
ภาพที่ 18 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 1



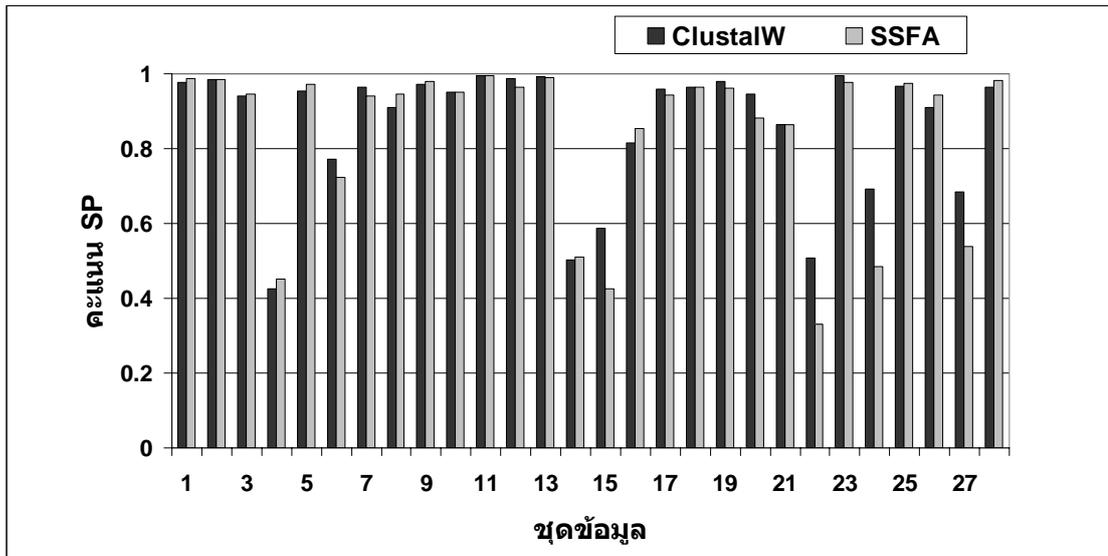
ภาพที่ 19 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 1



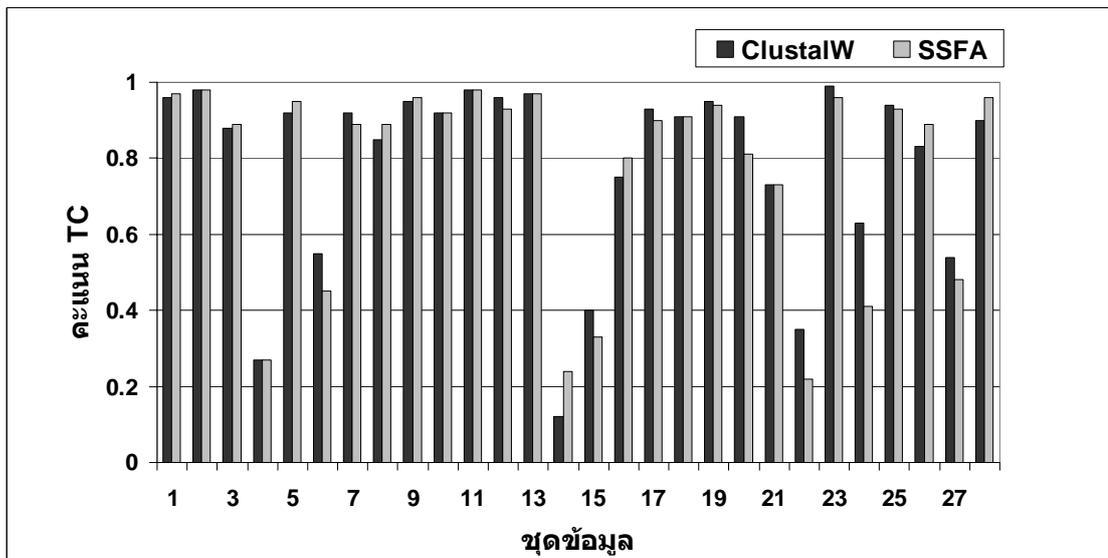
ภาพที่ 20 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 2



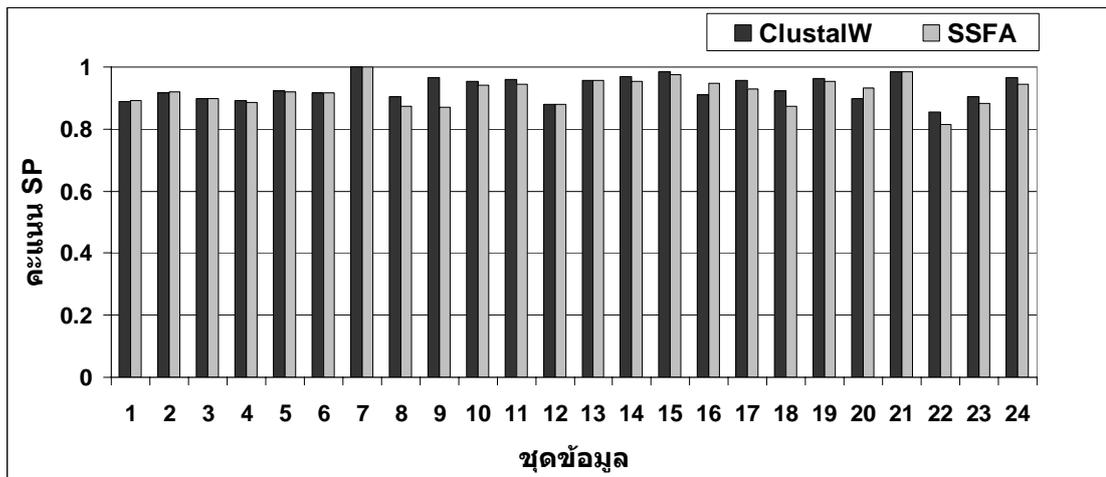
ภาพที่ 21 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 2



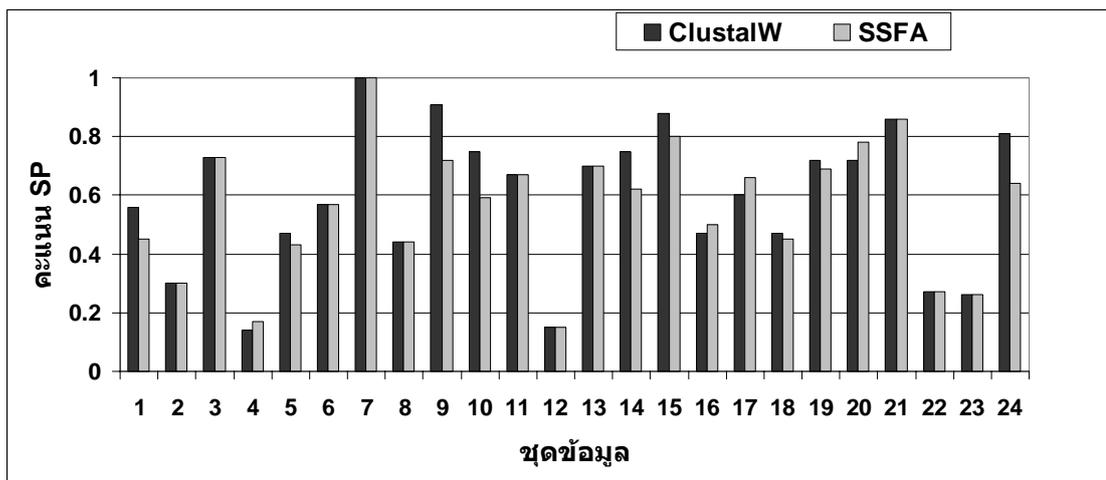
ภาพที่ 22 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 3



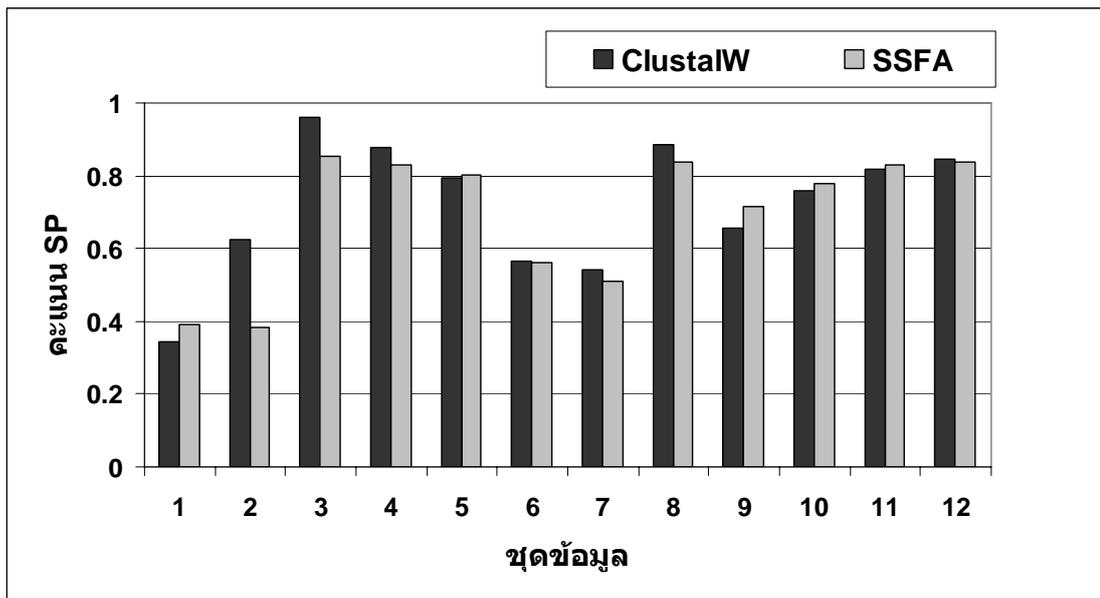
ภาพที่ 23 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 1 ชุดทดสอบ 3



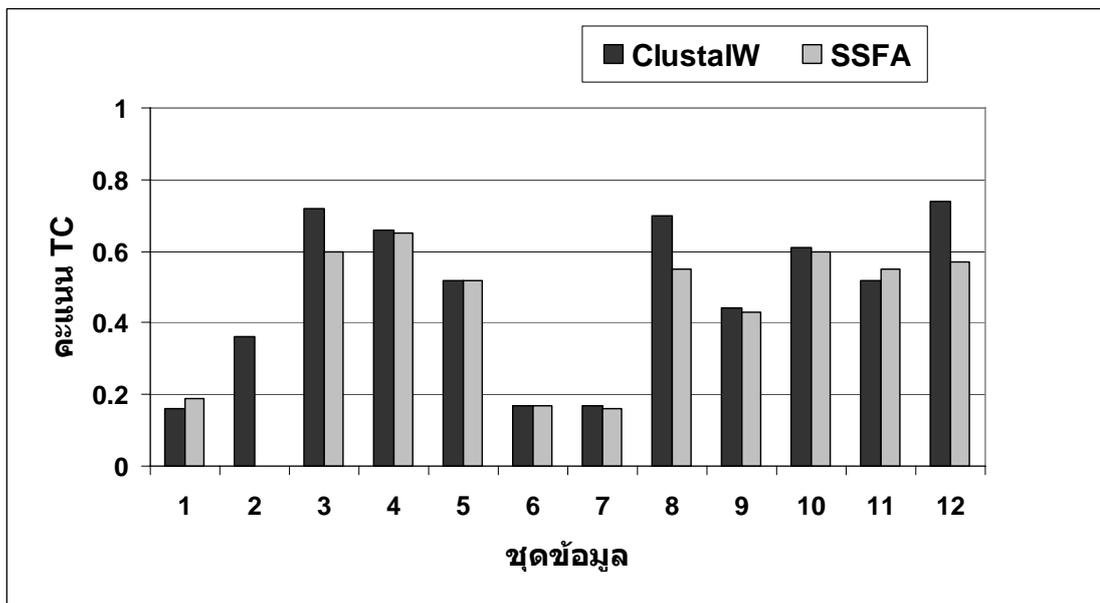
ภาพที่ 24 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 2



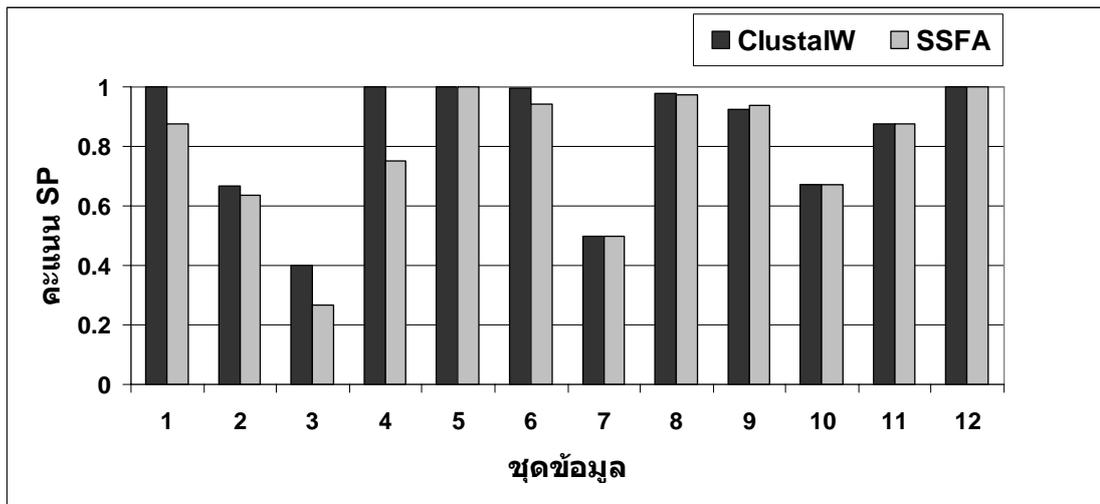
ภาพที่ 25 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 2



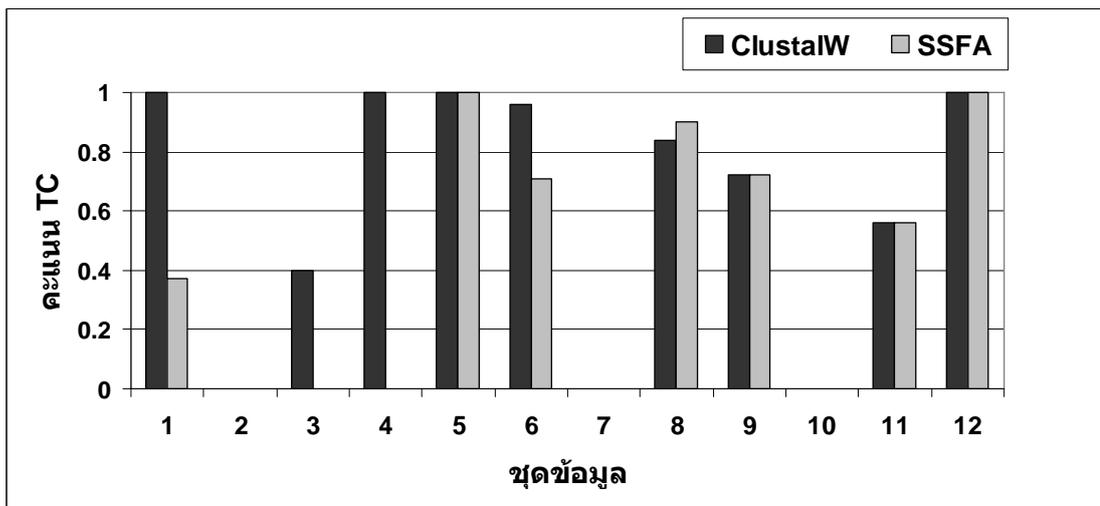
ภาพที่ 26 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 3



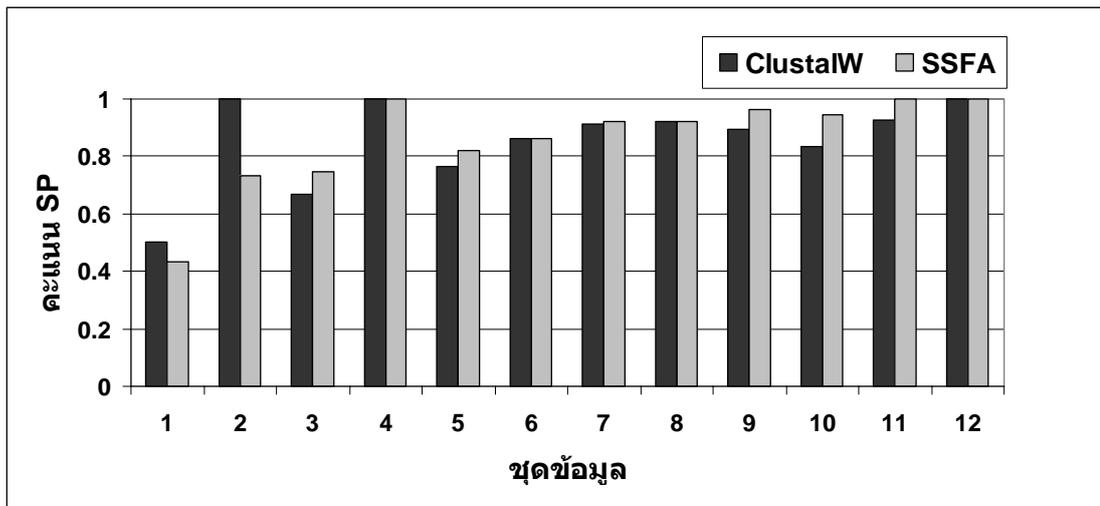
ภาพที่ 27 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 3



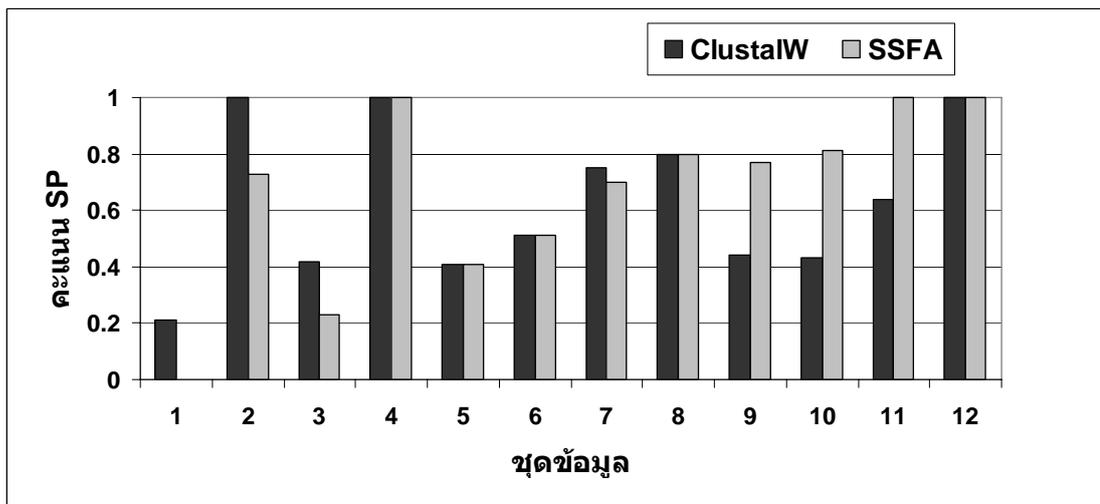
ภาพที่ 28 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 4



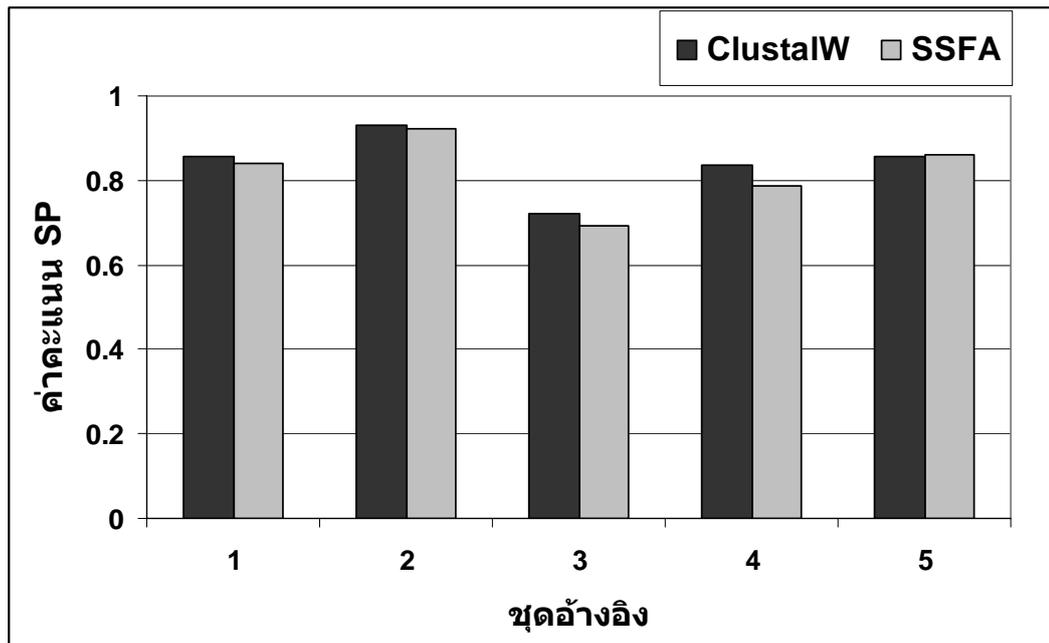
ภาพที่ 29 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 4



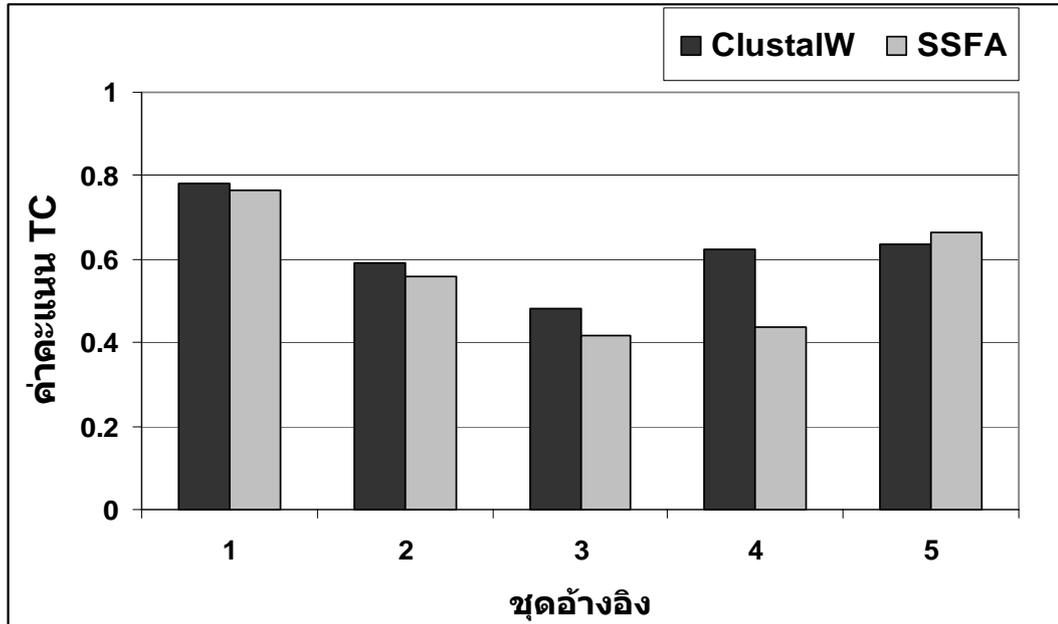
ภาพที่ 30 ค่าคะแนน SP ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 5



ภาพที่ 31 ค่าคะแนน TC ในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA กับชุดอ้างอิง: 5



ภาพที่ 32 ค่าคะแนน SP เฉลี่ยในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA ของแต่ละชุดอ้างอิง

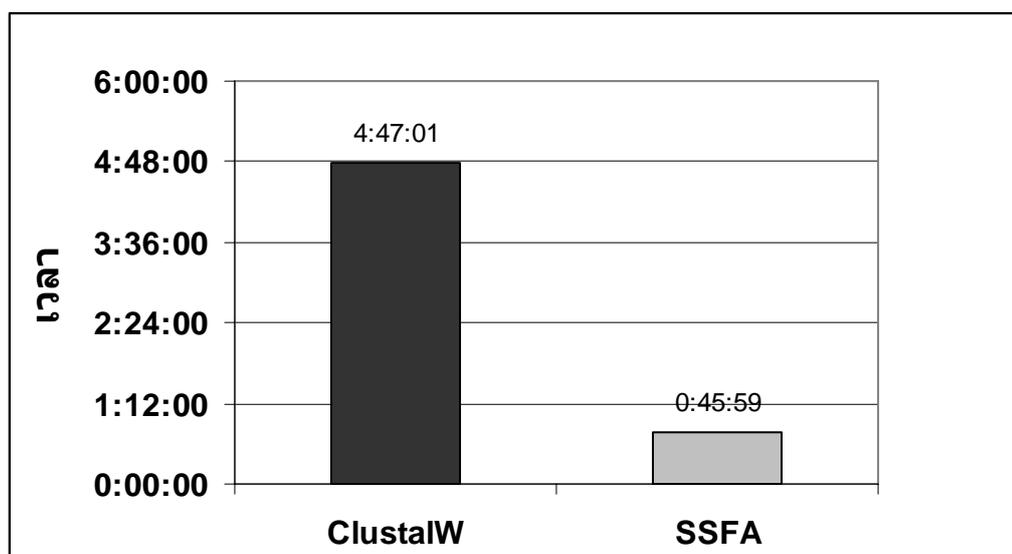


ภาพที่ 33 ค่าคะแนน TC เฉลี่ยในการเทียบเรียงกลุ่มลำดับข้อมูลของโปรแกรม ClustalW และโปรแกรม SSFA ของแต่ละชุดอ้างอิง

2. ผลการเทียบเรียงชุดข้อมูลของฐานข้อมูล PREFAB

การเทียบเรียงชุดข้อมูลของฐานข้อมูล PREFAB เปรียบเทียบประสิทธิภาพทางด้านเวลา ด้วยเวลารวมในการเทียบเรียงทุกชุดของฐานข้อมูลเช่นเดียวกับฐานข้อมูล BALiBASE ส่วนการเปรียบเทียบด้านความถูกต้องพิจารณาจากค่าคะแนน Q ทั้งจากการเทียบเรียงแต่ละชุดข้อมูล และค่าคะแนนเฉลี่ยของทุกชุดข้อมูล

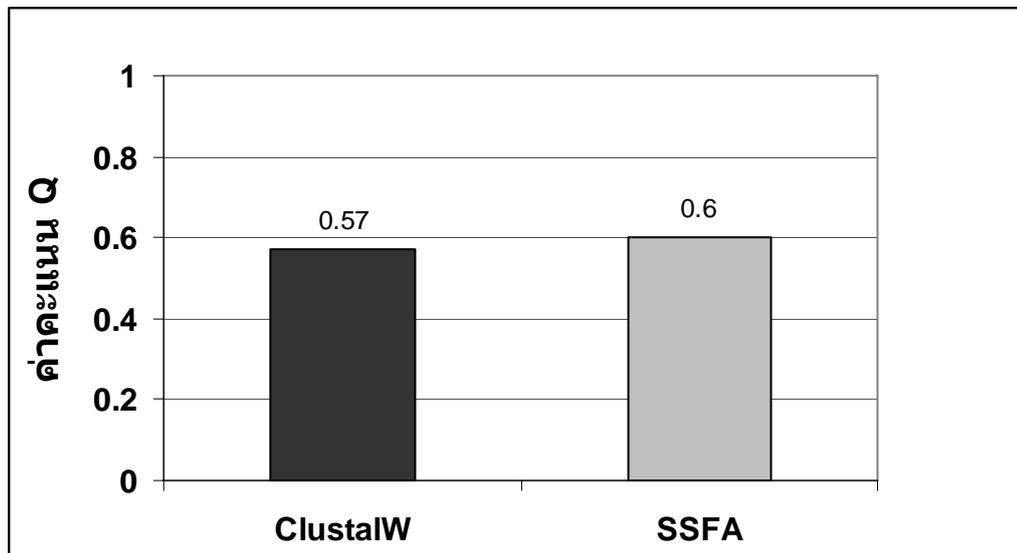
ภาพที่ 34 แสดงให้เห็นถึงความแตกต่างของเวลาที่โปรแกรม ClustalW และ โปรแกรม SSFA ใช้ในการเทียบเรียงชุดข้อมูลทั้ง 1682 ชุดของฐานข้อมูล PREFAB จะเห็นได้ว่าโปรแกรม ClustalW ใช้เวลาในการเทียบเรียงมากกว่าโปรแกรม SSFA มากถึง 6.24 เท่า ซึ่งความแตกต่างของเวลาที่ใช้ในการเทียบเรียงฐานข้อมูล PREFAB เห็นได้ชัดกว่าฐานข้อมูล BALiBASE เนื่องจากฐานข้อมูล PREFAB มีจำนวนชุดข้อมูลมากกว่า และในแต่ละชุดข้อมูลมีจำนวนลำดับข้อมูลมากกว่าฐานข้อมูล BALiBASE



ภาพที่ 34 เวลาที่โปรแกรม ClustalW และ โปรแกรม SSFA ใช้ในการเทียบเรียงฐานข้อมูล PREFAB

ในด้านความถูกต้องเมื่อพิจารณาจากค่าคะแนน Q เฉลี่ยของผลลัพธ์ที่ได้จากโปรแกรม ClustalW และ โปรแกรม SSFA ดังภาพที่ 35 ซึ่งมีความใกล้เคียงกัน โดยที่โปรแกรม SSFA ได้ค่าคะแนน Q เฉลี่ย 0.60 โดยมีค่าความเบี่ยงเบนมาตรฐานอยู่ที่ 0.34 ส่วนโปรแกรม ClustalW ได้ค่า

คะแนน Q เฉลี่ย 0.57 โดยมีค่าความเบี่ยงเบนมาตรฐานอยู่ที่ 0.33 ซึ่งถือว่ามีความใกล้เคียง และมีแนวโน้มของค่าคะแนนที่คล้ายกัน แสดงให้เห็นว่าโปรแกรม SSFA ให้ความถูกต้องในการเทียบเรียงกลุ่มลำดับข้อมูลในระดับเดียวกับโปรแกรม ClustalW

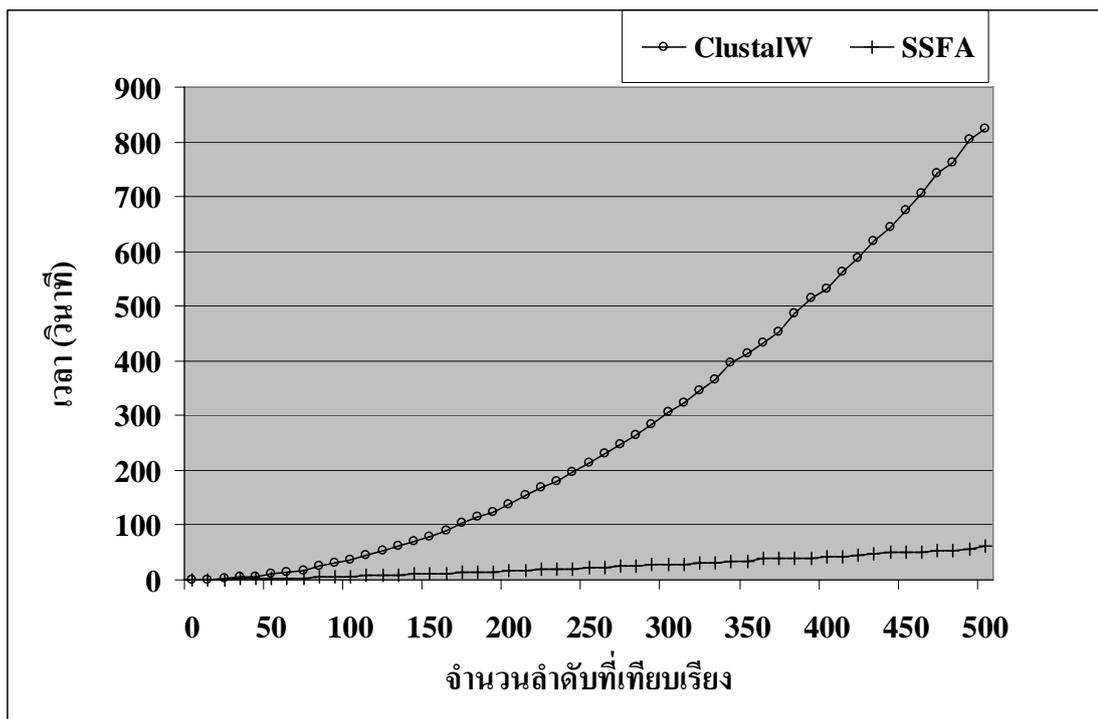


ภาพที่ 35 ค่าคะแนน Q เฉลี่ยในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA

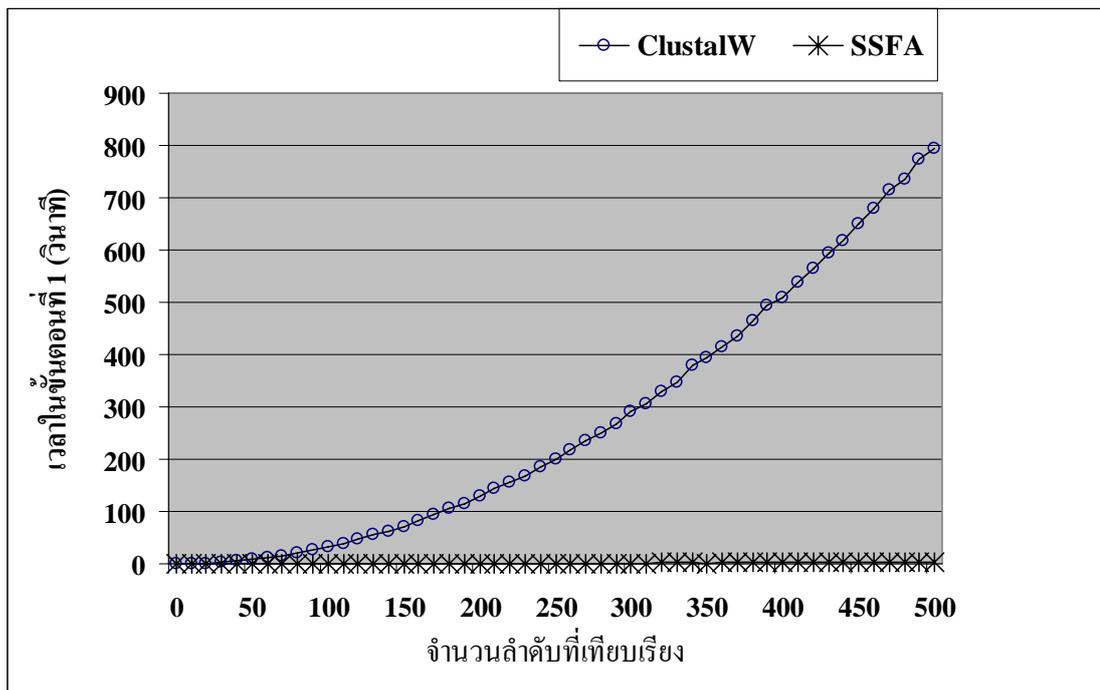
3. แนวโน้มทางด้านเวลาในการเทียบเรียงกลุ่มลำดับข้อมูลขนาดต่าง ๆ

ปัจจัยหลักที่ทำให้เวลาในการเทียบเรียงกลุ่มลำดับข้อมูลแตกต่างกัน คือจำนวนข้อมูลที่เข้าเทียบเรียง เมื่อจำนวนข้อมูลในกลุ่มลำดับข้อมูลมากขึ้นเวลาที่ใช้ในการเทียบเรียงจะมากขึ้นอย่างเห็นได้ชัด ในงานวิจัยนี้ได้ใช้ตัวอย่างข้อมูลโปรตีนจากฐานข้อมูล NCBI จัดกลุ่มลำดับที่เข้าเทียบเรียงให้มีจำนวนลำดับตั้งแต่ 10 – 500 ลำดับ โดยมีความยาวของลำดับเฉลี่ย 280 ตัวอักษร เมื่อนำมาเทียบเรียงด้วยโปรแกรม ClustalW และโปรแกรม SSFA เวลาที่ใช้ในการเทียบเรียงเป็นดังภาพที่ 36 และเวลาที่ใช้ในการทำงานขั้นตอนที่ 1 ของการเทียบเรียงข้อมูลเป็นดังภาพที่ 37 ซึ่งแสดงให้เห็นว่าเมื่อจำนวนข้อมูลที่เข้าเทียบเรียงมากขึ้นโปรแกรม ClustalW ใช้เวลาเพิ่มมากขึ้นมากเมื่อจำนวนลำดับที่เทียบเรียงเพิ่มมากขึ้น เนื่องจากในขั้นตอนการสร้างตารางค่าความแตกต่างใช้วิธีเทียบเรียงแต่ละคู่ลำดับข้อมูลซึ่งมีความซับซ้อนของเวลา $O(N^2L^2)$ เมื่อ N คือจำนวนลำดับที่เทียบเรียง และ L คือความยาวของลำดับที่เทียบเรียง ซึ่งแตกต่างกับโปรแกรม SSFA ที่ใช้การพิจารณาความถี่ส่วนย่อยของข้อมูลในการสร้างตารางค่าความแตกต่างซึ่งมีความซับซ้อนของเวลา $O(NL)$ ทำ

ให้อัตราการเพิ่มขึ้นของเวลาที่ใช้ต่อจำนวนของลำดับนั้นน้อย และเวลาที่ใช้ในการประมวลผลน้อยกว่าโปรแกรม ClustalW มาก โดยในการเทียบเรียงกลุ่มลำดับข้อมูลจำนวน 500 ลำดับที่มีความยาวเฉลี่ย 280 ตัวอักษร โปรแกรม SSFA สามารถสร้างตารางค่าความแตกต่างได้เร็วกว่าโปรแกรม ClustalW 198.75 เท่า และเทียบเรียงได้เร็วกว่าโปรแกรม ClustalW 13.14 เท่าซึ่งแสดงให้เห็นว่าโปรแกรม SSFA สามารถรองรับการเทียบเรียงกลุ่มลำดับข้อมูลที่มีจำนวนมากได้แม้ทำงานบนเครื่องคอมพิวเตอร์ทั่วไป



ภาพที่ 36 เวลาที่ใช้ในการเทียบเรียงของโปรแกรม ClustalW และโปรแกรม SSFA ในจำนวนลำดับที่ต่างกัน



ภาพที่ 37 เวลาที่ใช้ในขั้นตอนที่ 1 ของการเทียบเรียงของโปรแกรม ClustalW และ โปรแกรม SSFA ในจำนวนลำดับที่ต่างกัน

สรุปและข้อเสนอแนะ

สรุป

การเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพในปัจจุบันนอกจากคำนึงถึงความถูกต้องของการเทียบเรียงแล้ว ยังต้องคำนึงถึงเวลาในการเทียบเรียงกลุ่มลำดับข้อมูลซึ่งแปรผันตรงกับจำนวนข้อมูลที่เพิ่มขึ้นอย่างรวดเร็วในปัจจุบัน งานวิจัยนี้จึงนำเสนอวิธีการเทียบเรียงกลุ่มลำดับข้อมูลแบบก้าวหน้าโดยพิจารณาความถี่ส่วนย่อยของกลุ่มลำดับข้อมูล ซึ่งงานวิจัยนี้ใช้เวลาในการประมวลผลน้อยมากโดยที่ความถูกต้องสามารถเทียบเคียงได้กับวิธีการเทียบเรียงกลุ่มลำดับข้อมูลที่นิยมในปัจจุบัน งานวิจัยนี้ใช้เทคนิค n-gram แบ่งส่วนย่อยของข้อมูลเพื่อหาความถี่ของแต่ละส่วนย่อยของกลุ่มลำดับข้อมูล สำหรับนำไปสร้างตารางค่าความแตกต่าง เพื่อใช้อ้างอิงในการเทียบเรียงกลุ่มลำดับข้อมูลแบบก้าวหน้า ซึ่งเทคนิคนี้ใช้เวลาในการทำงานน้อย ทำให้สามารถรองรับการเทียบเรียงกลุ่มลำดับข้อมูลจำนวนมาก และลำดับข้อมูลที่มีขนาดยาวได้

ผลลัพธ์การเทียบเรียงของงานวิจัยนี้เปรียบเทียบกับโปรแกรม ClustalW ซึ่งเป็นโปรแกรมในกลุ่มวิธีเทียบเรียงกลุ่มลำดับข้อมูลแบบก้าวหน้าที่เป็นที่นิยมที่สุดในปัจจุบัน ผลของการเทียบเรียงฐานข้อมูล BALiBASE โปรแกรม SSFA สามารถเทียบเรียงได้เร็วกว่าโปรแกรม ClustalW 2.72 เท่า โดยที่ความถูกต้องของการเทียบเรียงใกล้เคียงกันอยู่ในระดับที่ยอมรับได้ โปรแกรม SSFA ได้เปรียบโปรแกรม ClustalW ในข้อมูลชุดที่ 5 ที่เป็นกลุ่มข้อมูลที่มีลักษณะการแทรกภายใน และโปรแกรม SSFA จะเสียเปรียบโปรแกรม ClustalW ในข้อมูลชุดที่ 4 ซึ่งมีข้อมูลขนาดใหญ่เพิ่มเติมที่ตอนต้นหรือตอนท้ายของข้อมูล ซึ่งข้อมูลลักษณะนี้มีผลทำให้ส่วนจำเพาะที่โปรแกรม SSFA หาได้กระจายไปในส่วนข้อมูลที่เพิ่มเติมขึ้นมาทำให้ผลที่ได้แตกต่างจากผลเทียบเรียงอ้างอิง

สำหรับผลการเทียบเรียงเมื่อเทียบเรียงฐานข้อมูล PREFAB โปรแกรม SSFA สามารถเทียบเรียงได้เร็วกว่าโปรแกรม ClustalW 6.24 เท่า ซึ่งมากกว่าการเทียบเรียงฐานข้อมูล BALiBASE เนื่องจากฐานข้อมูล PREFAB มีชุดข้อมูลที่มากกว่า และในแต่ละชุดข้อมูลมีจำนวนลำดับและความยาวของสายลำดับมากกว่าฐานข้อมูล BALiBASE ซึ่งแสดงให้เห็นว่าเมื่อกลุ่มข้อมูลที่เข้าเทียบเรียงมีขนาดใหญ่ขึ้นประสิทธิภาพทางด้านเวลาของโปรแกรม SSFA จะดีกว่าโปรแกรม ClustalW เพิ่มขึ้นด้วย ในส่วนความถูกต้องของผลการเทียบเรียงนั้นโปรแกรม SSFA ให้ค่าความถูกต้องเฉลี่ยและค่าเบี่ยงเบนมาตรฐานใกล้เคียงกับ โปรแกรม ClustalW แสดงให้เห็นว่าความถูกต้องของโปรแกรม SSFA นั้นอยู่ในระดับเดียวกับโปรแกรม ClustalW

เมื่อพิจารณาประสิทธิภาพทางด้านเวลาเมื่อเทียบเรียงกลุ่มข้อมูลขนาดต่างๆ แสดงให้เห็นว่าเมื่อกลุ่มข้อมูลที่เข้าเทียบเรียงมีขนาดใหญ่ขึ้นจะยิ่งเห็นความแตกต่างของเวลาที่ใช้เทียบเรียงมากขึ้น โดยเมื่อเทียบเรียงกลุ่มลำดับข้อมูลจำนวน 500 ลำดับ โปรแกรม SSFA สามารถเทียบเรียงได้เร็วกว่าโปรแกรม ClustalW ถึง 13.14 เท่า ซึ่งแสดงให้เห็นว่าเทคนิคการเทียบเรียงกลุ่มลำดับข้อมูลโดยพิจารณาความถี่ส่วนย่อยของลำดับสามารถรองรับการเทียบเรียงกลุ่มข้อมูลขนาดใหญ่ได้ดีกว่าโปรแกรม ClustalW

ข้อเสนอแนะ

วิทยานิพนธ์ฉบับนี้นำเสนอวิธีการเทียบเรียงกลุ่มลำดับข้อมูลโดยพิจารณาความถี่ส่วนย่อยของลำดับ ซึ่งเป็นวิธีที่ใช้เวลาในการประมวลผลน้อย และรองรับการเทียบเรียงกลุ่มลำดับข้อมูลขนาดใหญ่ โดยที่มีความถูกต้องในระดับที่ยอมรับได้ อย่างไรก็ตามการเทียบเรียงกลุ่มลำดับข้อมูลยังคงต้องการความถูกต้องแม่นยำที่สูงกว่าปัจจุบัน ซึ่งแนวทางการพัฒนาในอนาคตควรเป็นการพัฒนาในขั้นตอนที่เหลือเพื่อเพิ่มความถูกต้องแม่นยำในการเทียบเรียงให้สูงขึ้น โดยที่ยังใช้เวลาประมวลผลน้อย และรองรับการเทียบเรียงข้อมูลขนาดใหญ่เช่นกัน

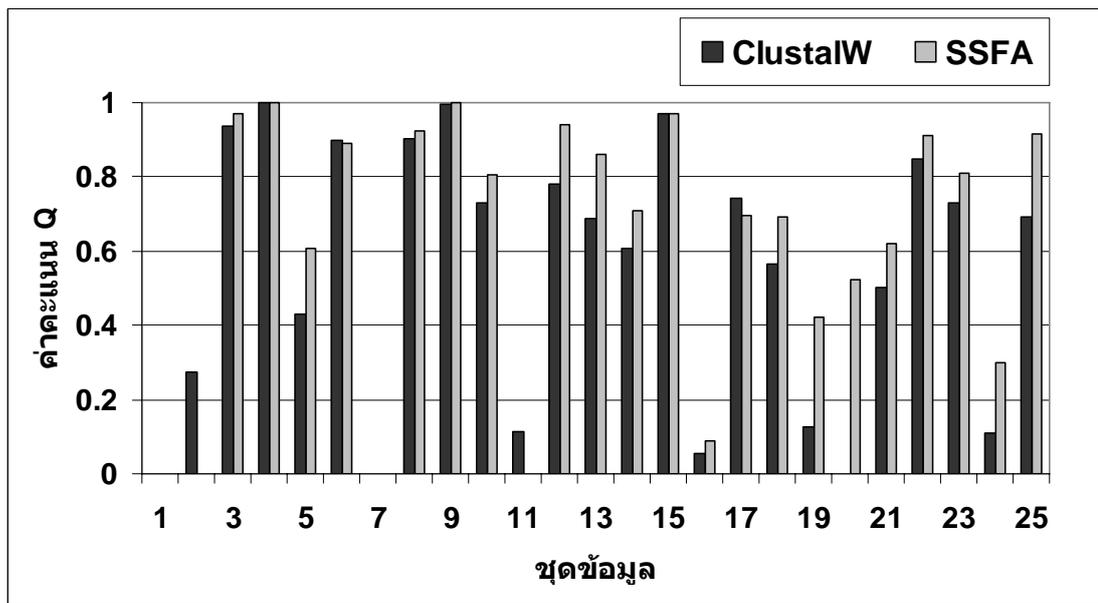
เอกสารและสิ่งอ้างอิง

- คมสัน จันมา และ พันธุ์ปิติ เปี่ยมสง่า. 2548. การเทียบเรียงกลุ่มลำดับข้อมูลทางชีวภาพโดยพิจารณาความถี่ส่วนย่อยของลำดับ, น. 231-240. ใน **National Computer Science and Engineering Conference ครั้งที่ 9**. มหาวิทยาลัยหอการค้าไทย, กรุงเทพฯ.
- จิตมณฑิ์ เขียนดวงจันทร์. 2544. ขั้นตอนวิธีในการจับคู่ลำดับเบสดีเอ็นเอ. วิทยานิพนธ์ปริญญาโท, มหาวิทยาลัยเกษตรศาสตร์.
- พจน์ ศรีบุญลือ, โสพิศ วงศ์คำ และ พัชรี บุญศิริ. 2543. ตำราชีวเคมี. ภาควิชาชีวเคมี คณะแพทยศาสตร์ มหาวิทยาลัยขอนแก่น, ขอนแก่น.
- ไพศาล เหล่าสุวรรณ. 2539. พันธุศาสตร์. ไทยวัฒนาพานิช, กรุงเทพฯ.
- มุกดา ณีภูธรสมบุญ. 2526. พันธุศาสตร์พื้นฐาน. ภาควิชาพฤกษศาสตร์ คณะวิทยาศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย, กรุงเทพฯ.
- วีรวุฒิ คงบุญเกียรติ. 2545. การเทียบเรียงกลุ่มลำดับข้อมูลชีวภาพโดยใช้ความถี่แบบพิจารณาค่าน้ำหนัก. วิทยานิพนธ์ปริญญาโท, มหาวิทยาลัยเกษตรศาสตร์.
- อมรา กัมภีรานนท์. 2540. พันธุศาสตร์ของเซลล์. มหาวิทยาลัยเกษตรศาสตร์, กรุงเทพฯ.
- Baldi, P. and S. Brunak. 1999. **Bioinformatics : the machine learning approach**. 3 ed. The MIT press, London, England.
- Crochemore, M., G.M. Landau and M. Ziv-Ukelson. 2003. A Subquadratic Sequence Alignment Algorithm for Unrestricted Scoring Matrices. **SIAM Journal on Computing** 32 (6): 1654-1673.

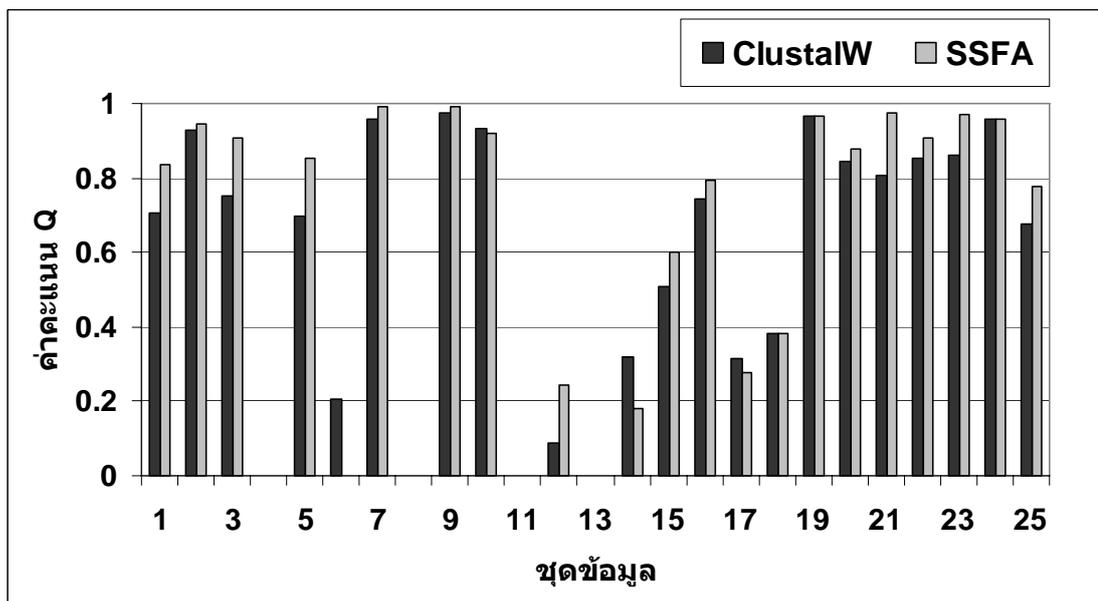
- Edgar, R.C. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. **Nucleic Acids Research** 32 (5): 1792-1797.
- Gusfield, D. 1999. **Algorithms on strings, trees, and sequences : computer science and computation biology..** 2 ed. The press syndicate of The University of Cambridge, New York, USA.
- Notredame, C. 2002. Recent progress in multiple sequence alignment: a survey. **Pharmacogenomics** 3 (1): 131-144.
- Saitou, N. and M. Nei. 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees.. **Molecular Biology and Evolution** 4 (4): 406-425.
- Thompson, J.D., D.G. Higgins and T.J. Gibson. 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice.. **Nucleic Acids Research** 22 (22): 4673-4680.
- _____, _____, _____, F. Plewniak and F. Jeanmougin. 1997. The CLUSTAL_X windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. **Nucleic Acids Research** 25 (24): 4876-4882
- _____, _____ and O. Poch. 1999. A comprehensive comparison of multiple sequence alignment programs. **Nucleic Acids Research** 27 (13): 2682-2690.

ภาคผนวก

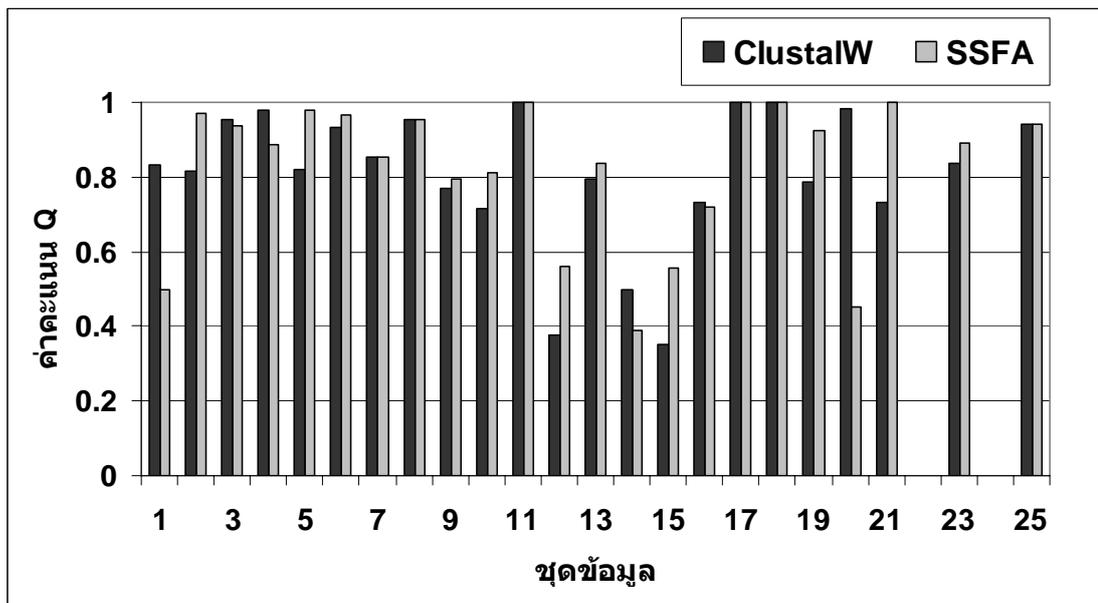
ผลการเทียบเรียงชุดข้อมูลของฐานข้อมูล PREFAB แยกตามกลุ่มข้อมูล



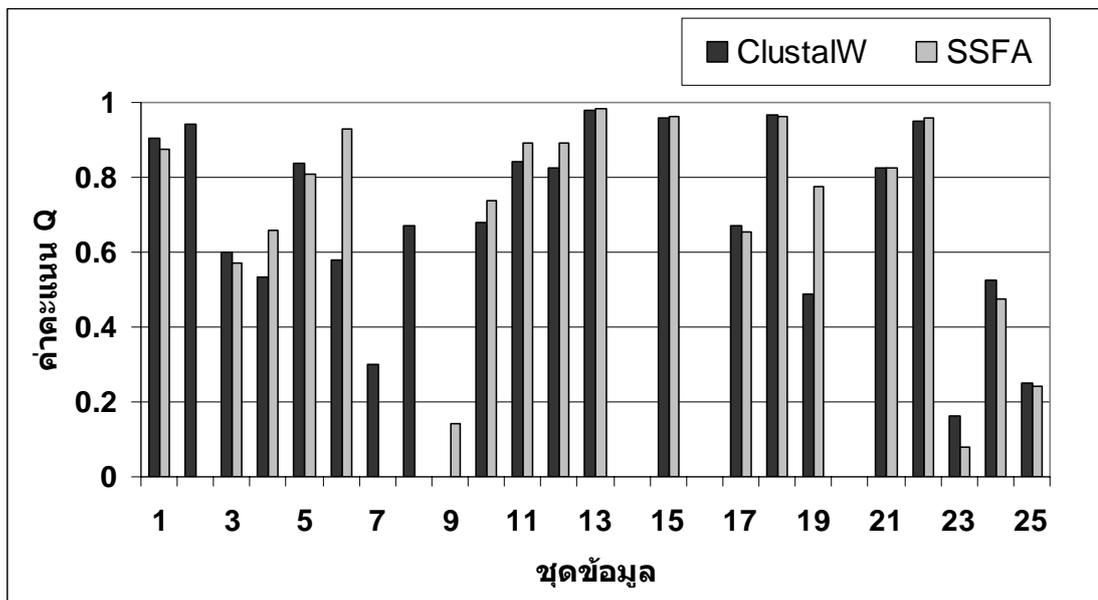
ภาพผนวกที่ 1 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 1



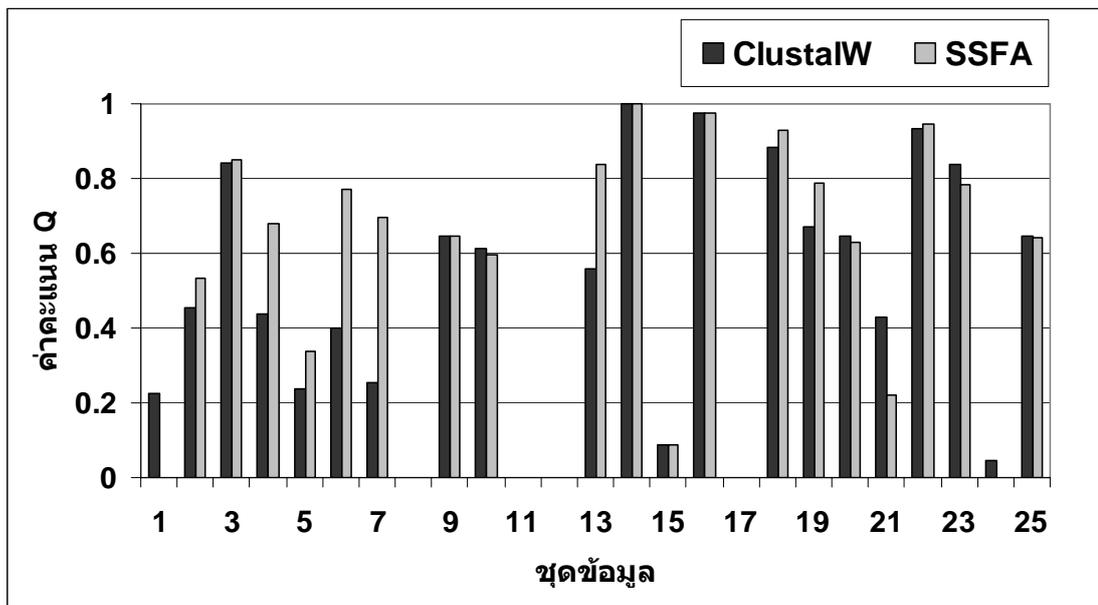
ภาพผนวกที่ 2 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 2



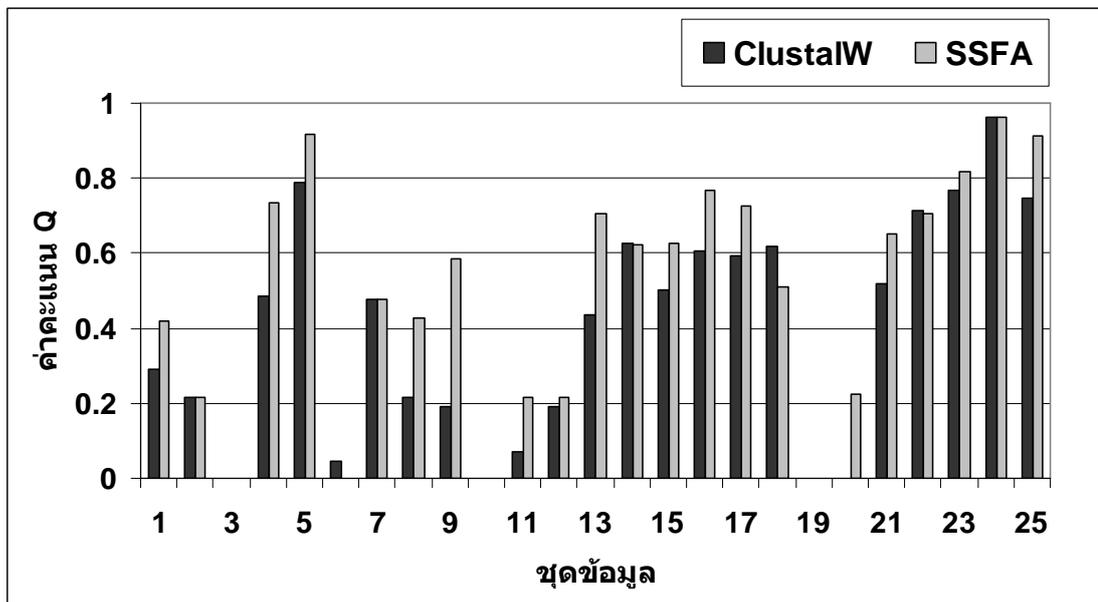
ภาพผนวกที่ 3 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 3



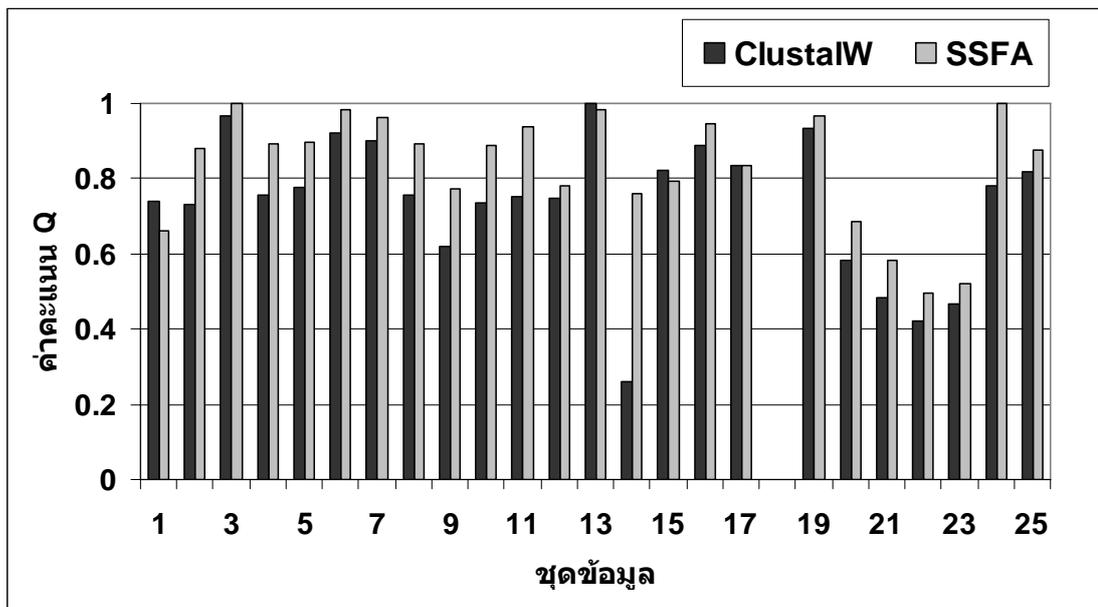
ภาพผนวกที่ 4 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 4



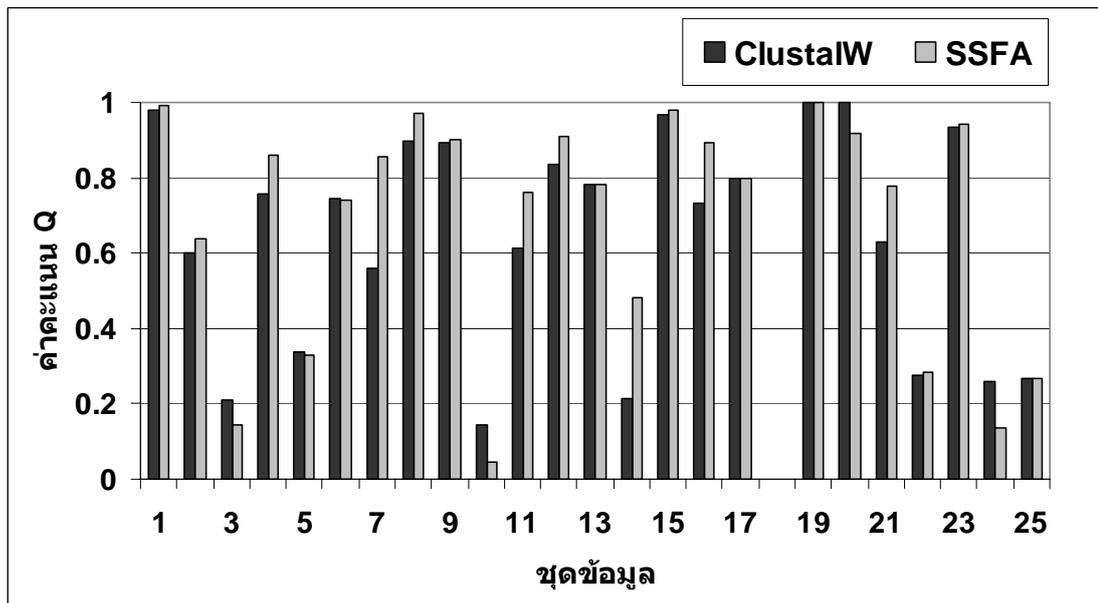
ภาพผนวกที่ 5 ค่าคะแนน Q ในการเทียบเคียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 5



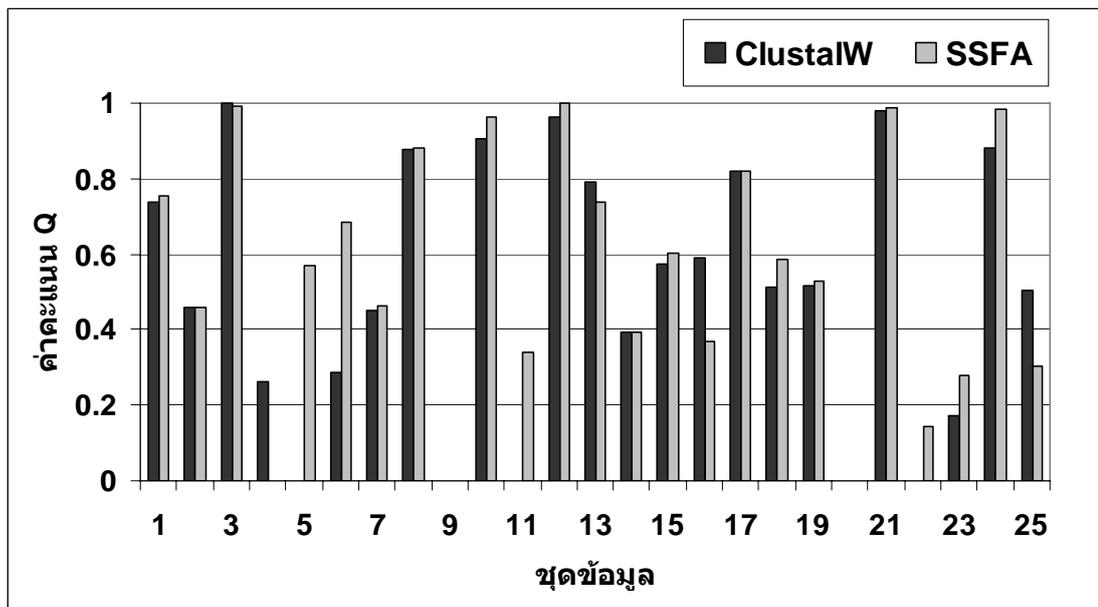
ภาพผนวกที่ 6 ค่าคะแนน Q ในการเทียบเคียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 6



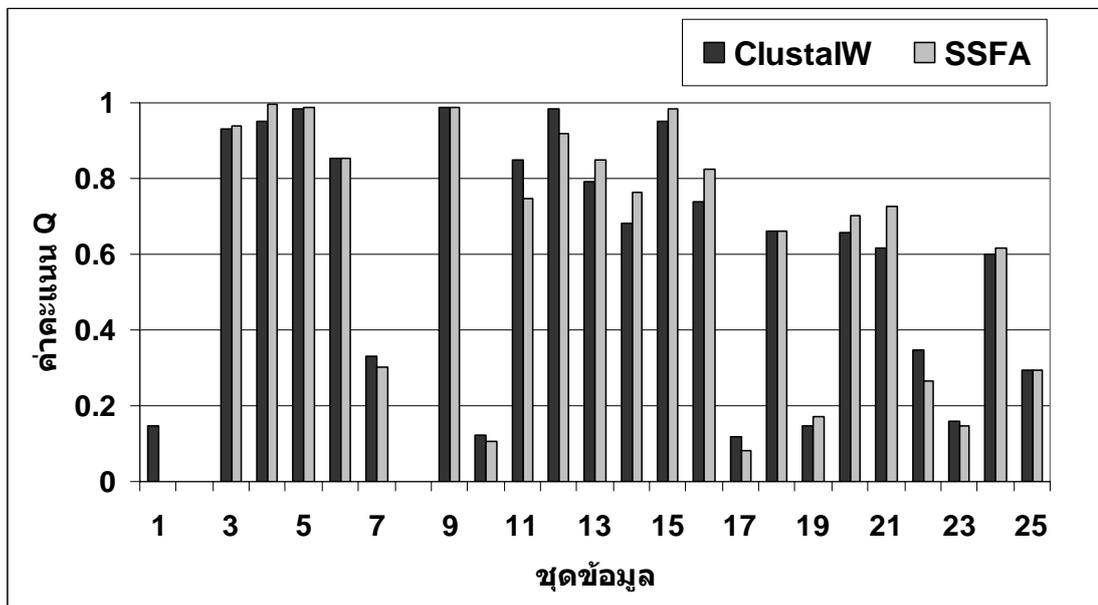
ภาพผนวกที่ 7 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 7



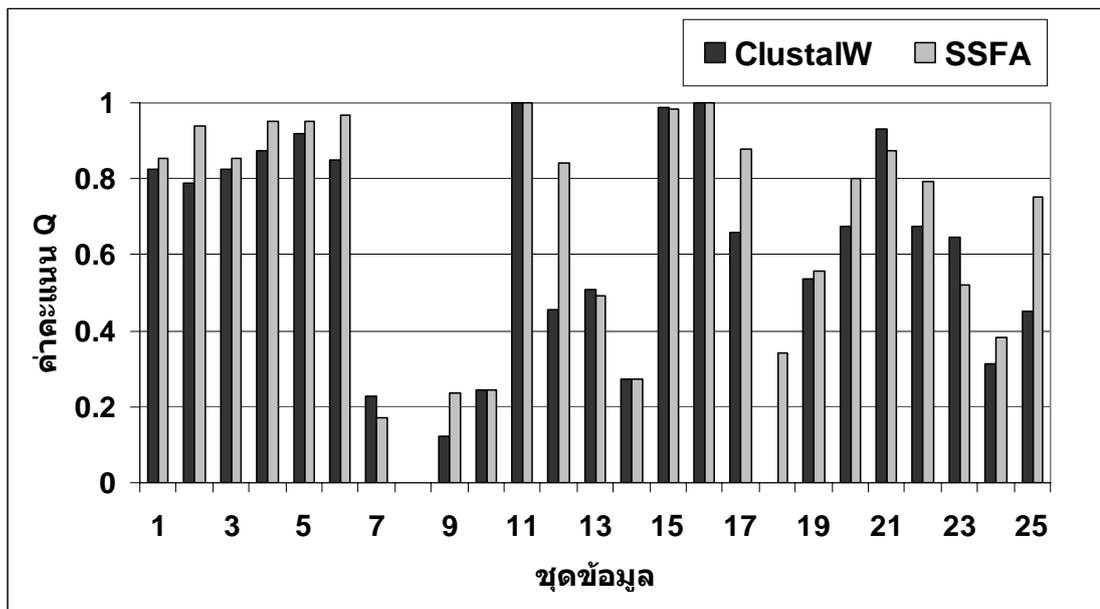
ภาพผนวกที่ 8 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 8



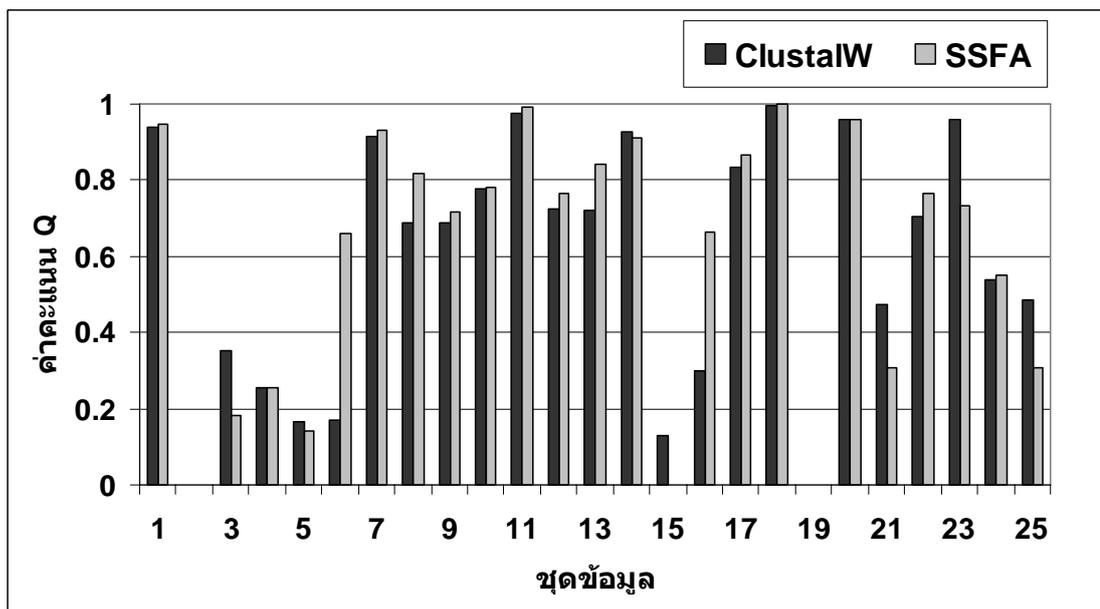
ภาพผนวกที่ 9 ค่าคะแนน Q ในการเทียบเคียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 9



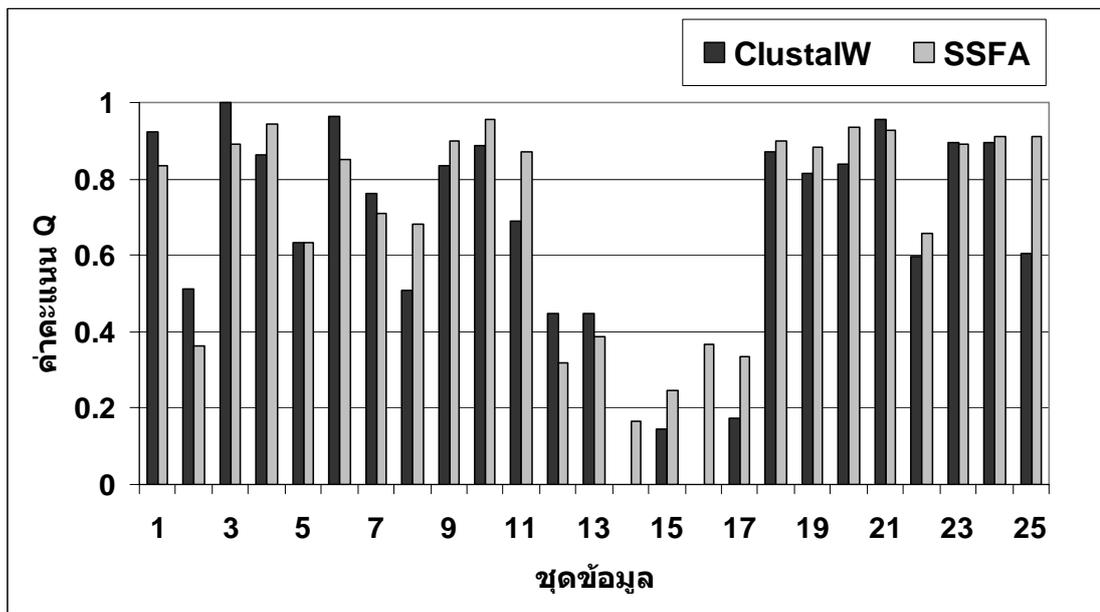
ภาพผนวกที่ 10 ค่าคะแนน Q ในการเทียบเคียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 10



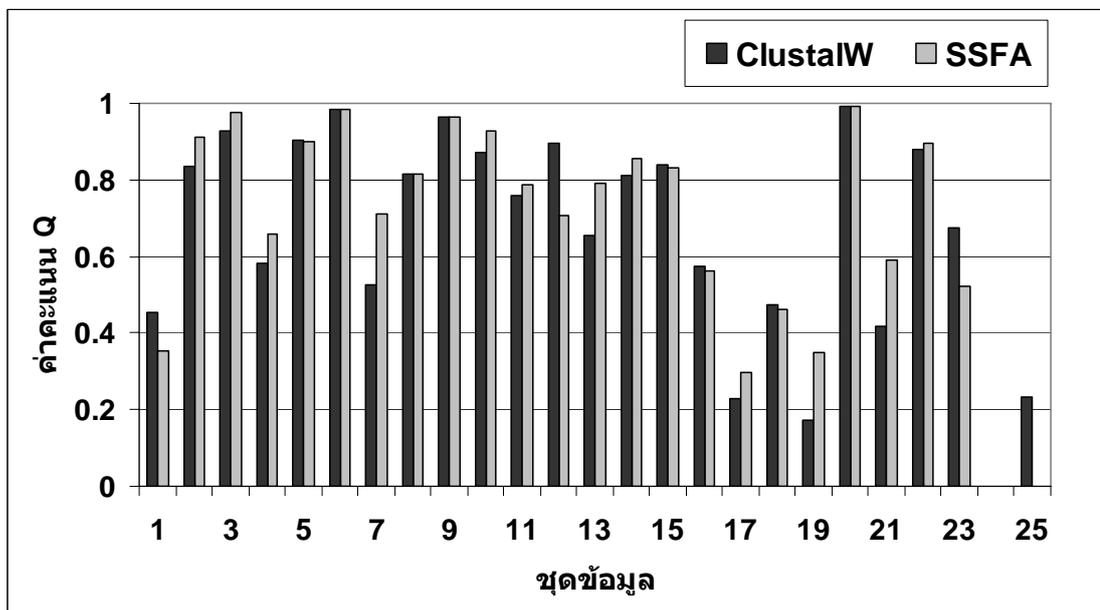
ภาพผนวกที่ 11 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 11



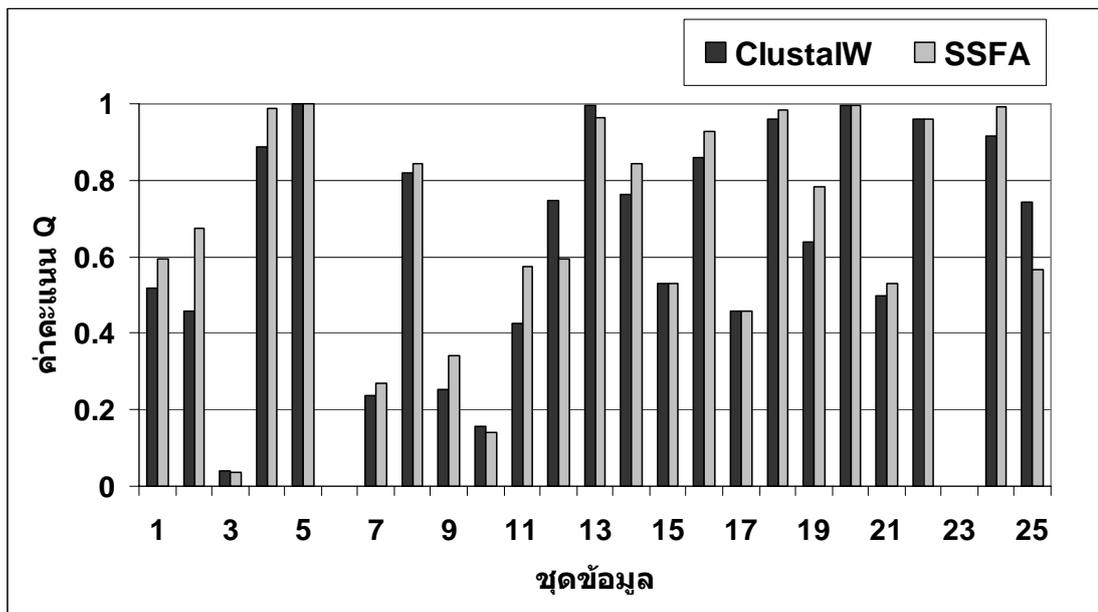
ภาพผนวกที่ 12 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 12



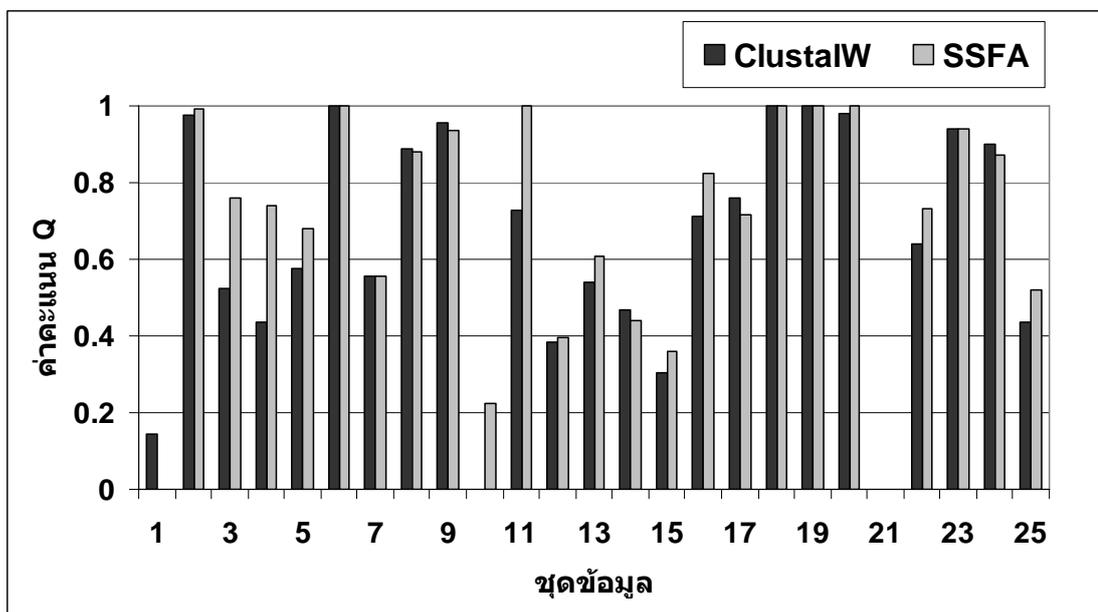
ภาพผนวกที่ 13 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 13



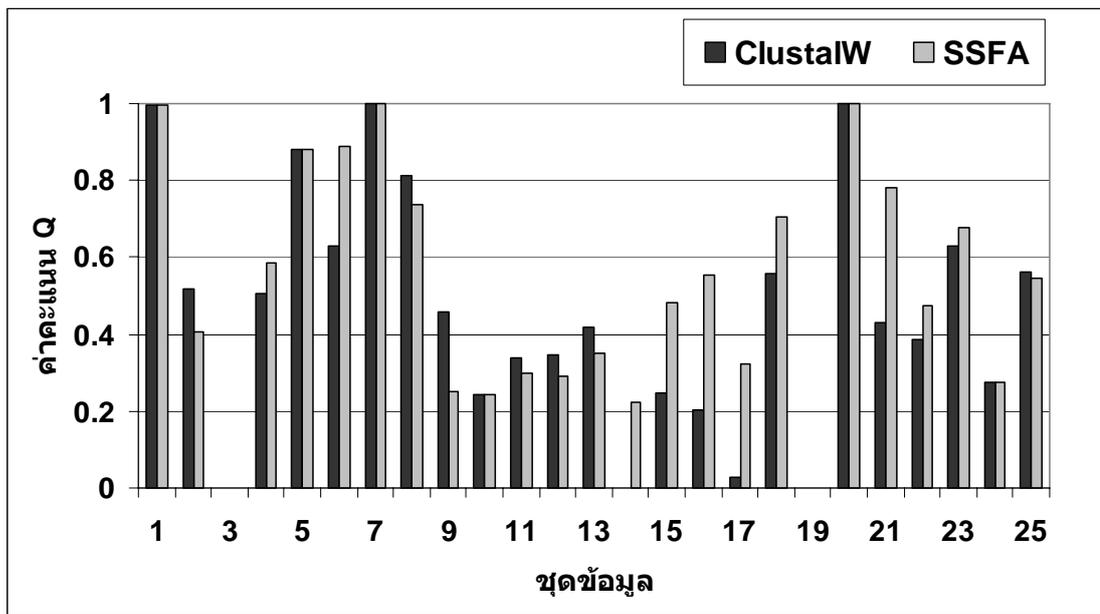
ภาพผนวกที่ 14 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 14



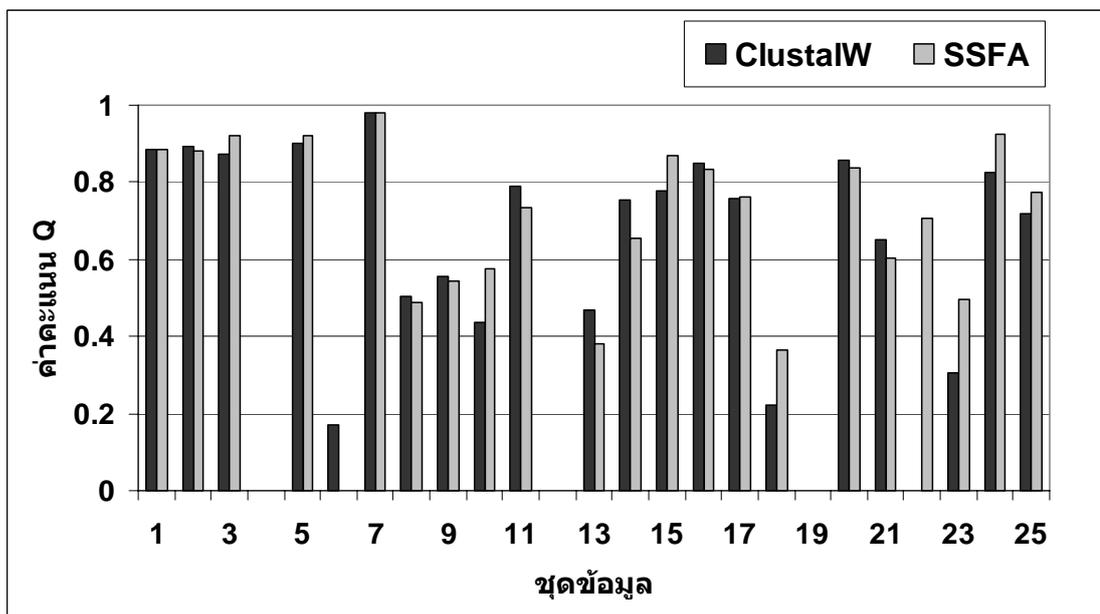
ภาพผนวกที่ 15 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 15



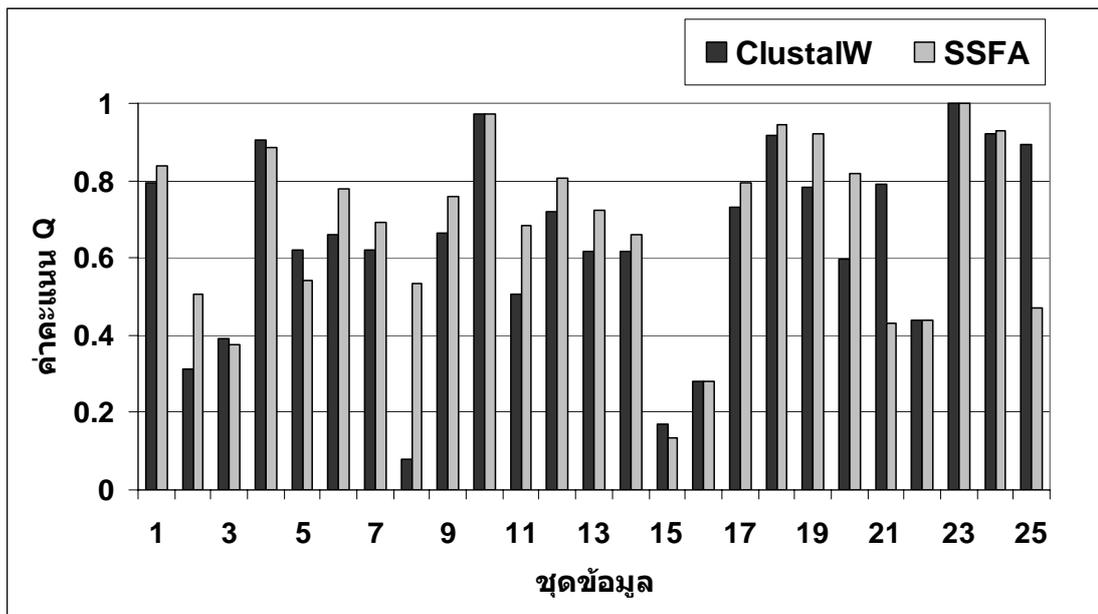
ภาพผนวกที่ 16 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 16



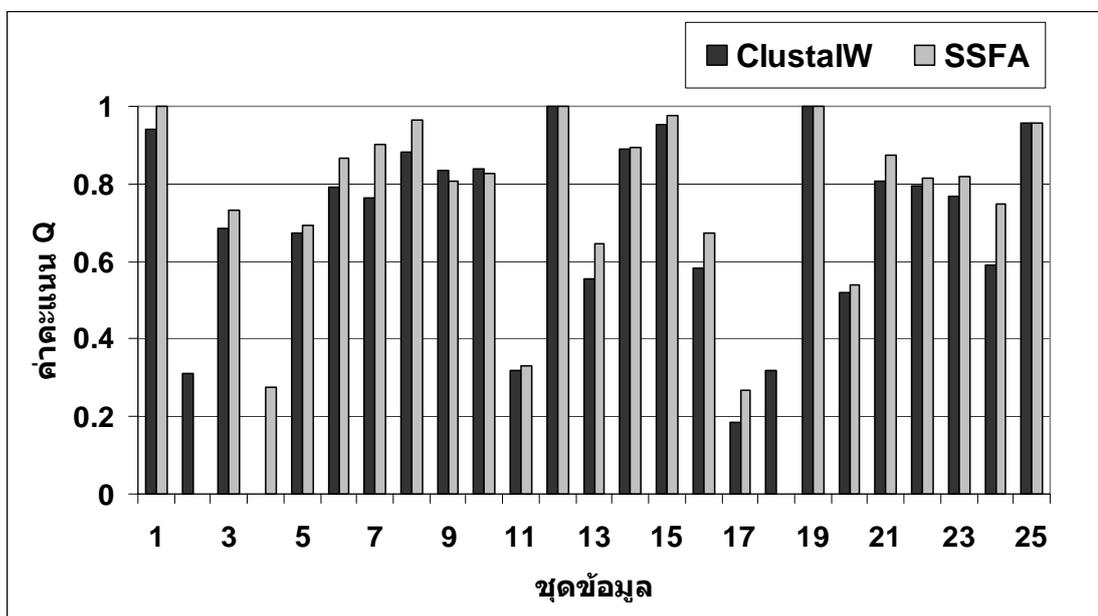
ภาพผนวกที่ 17 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 17



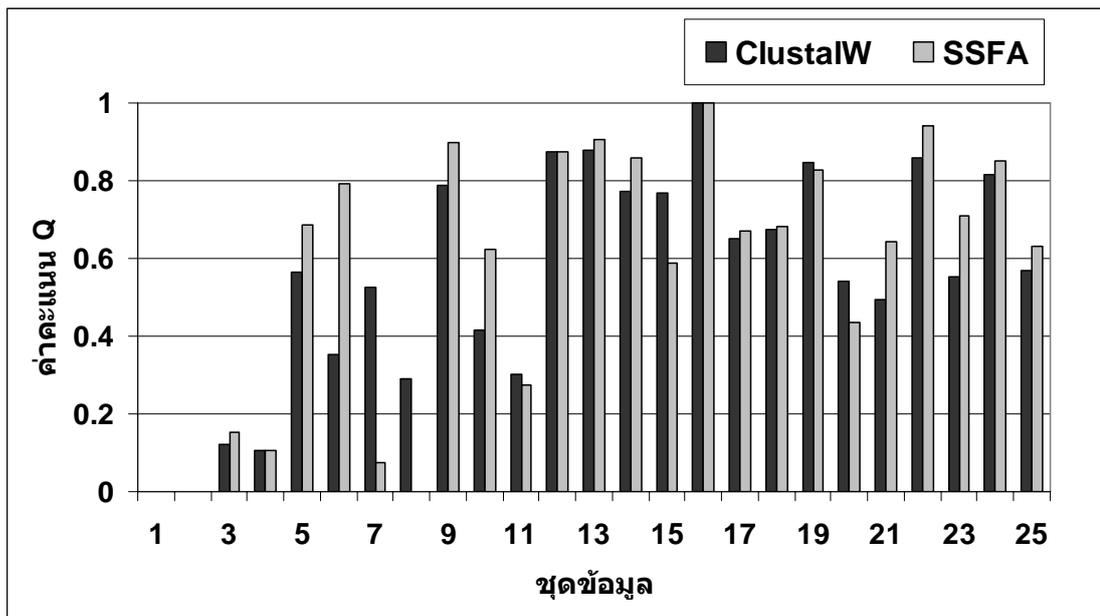
ภาพผนวกที่ 18 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 18



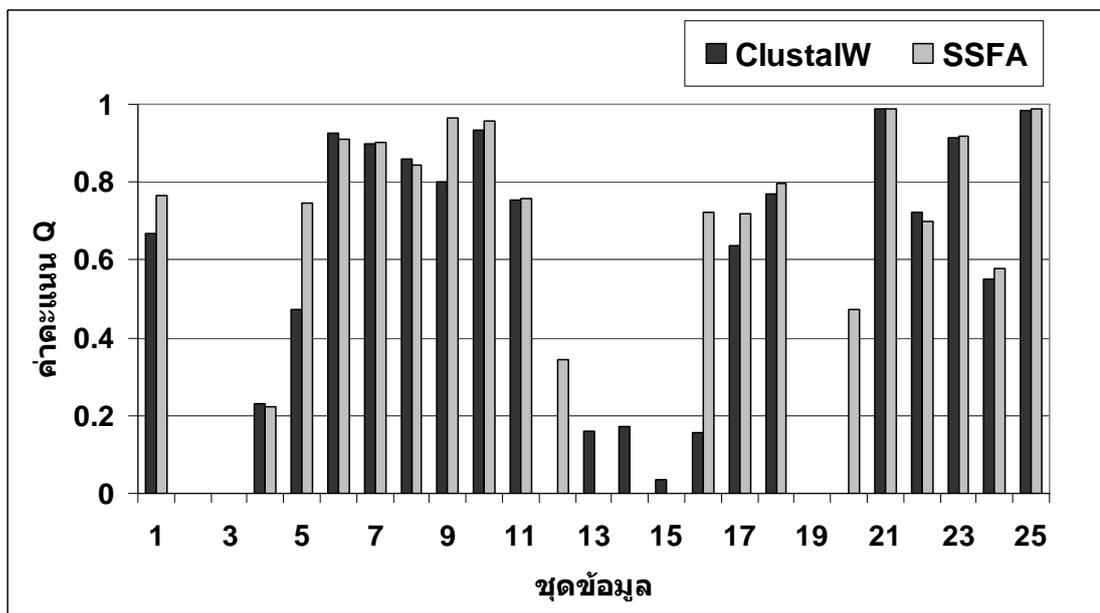
ภาพผนวกที่ 19 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 19



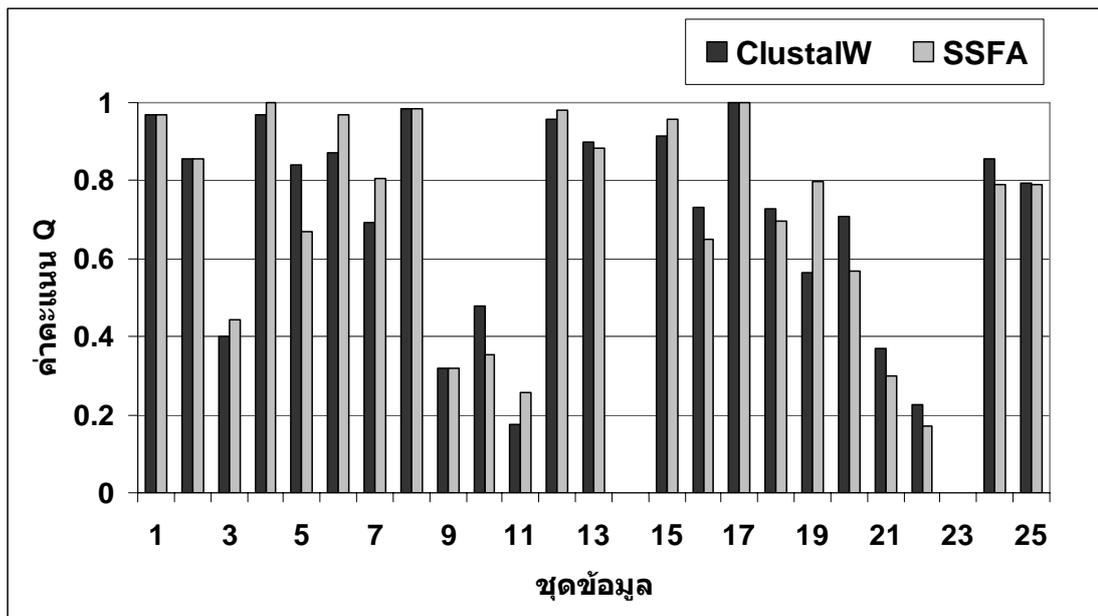
ภาพผนวกที่ 20 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 20



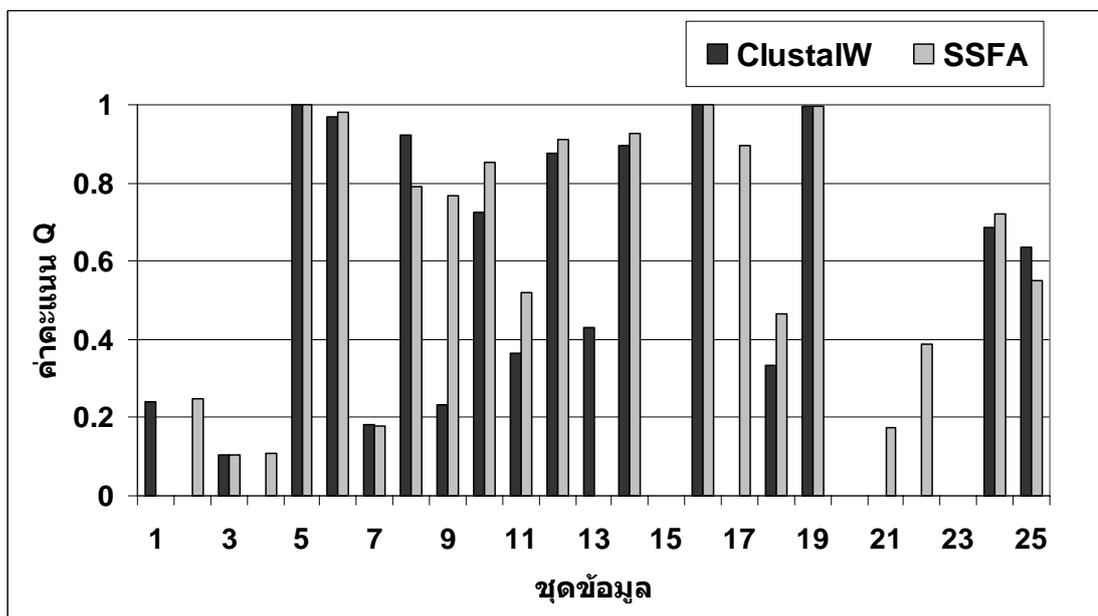
ภาพผนวกที่ 21 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 21



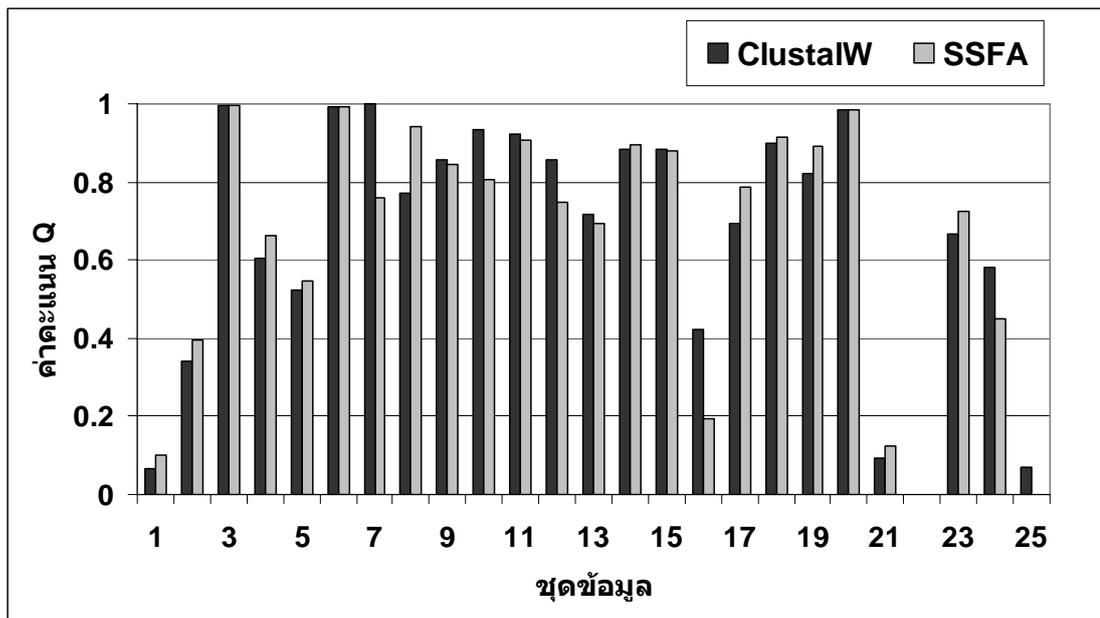
ภาพผนวกที่ 22 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 22



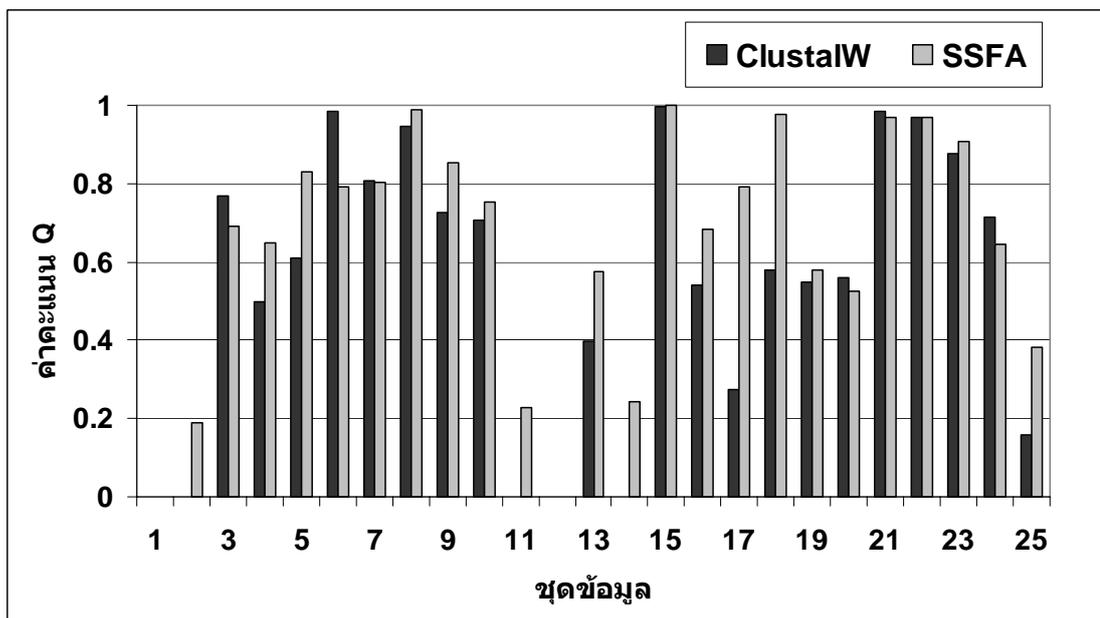
ภาพผนวกที่ 23 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 23



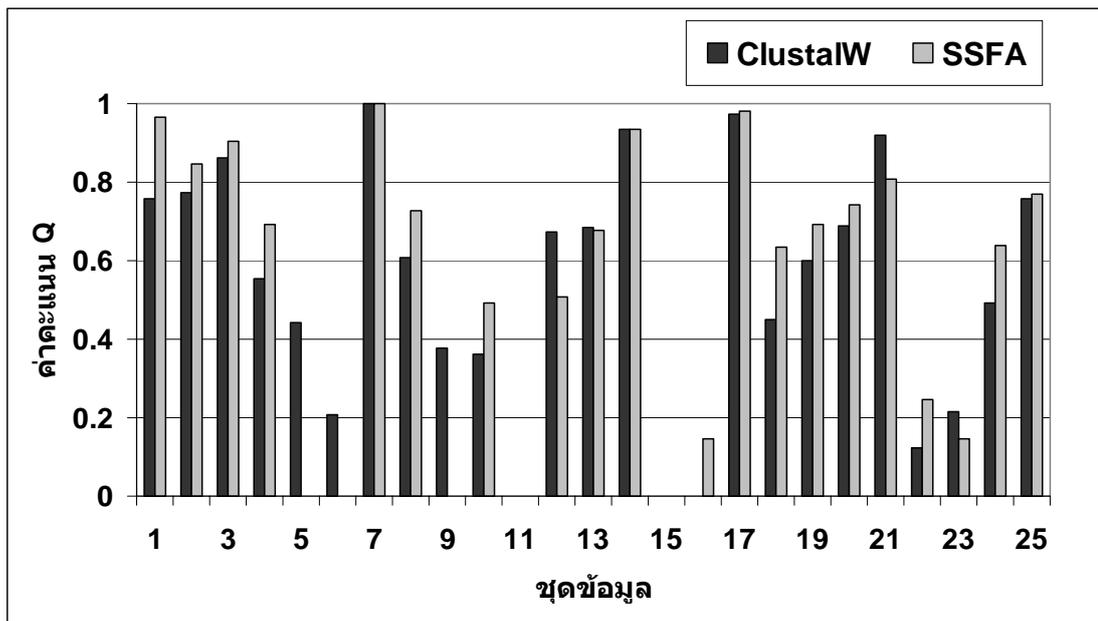
ภาพผนวกที่ 24 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 24



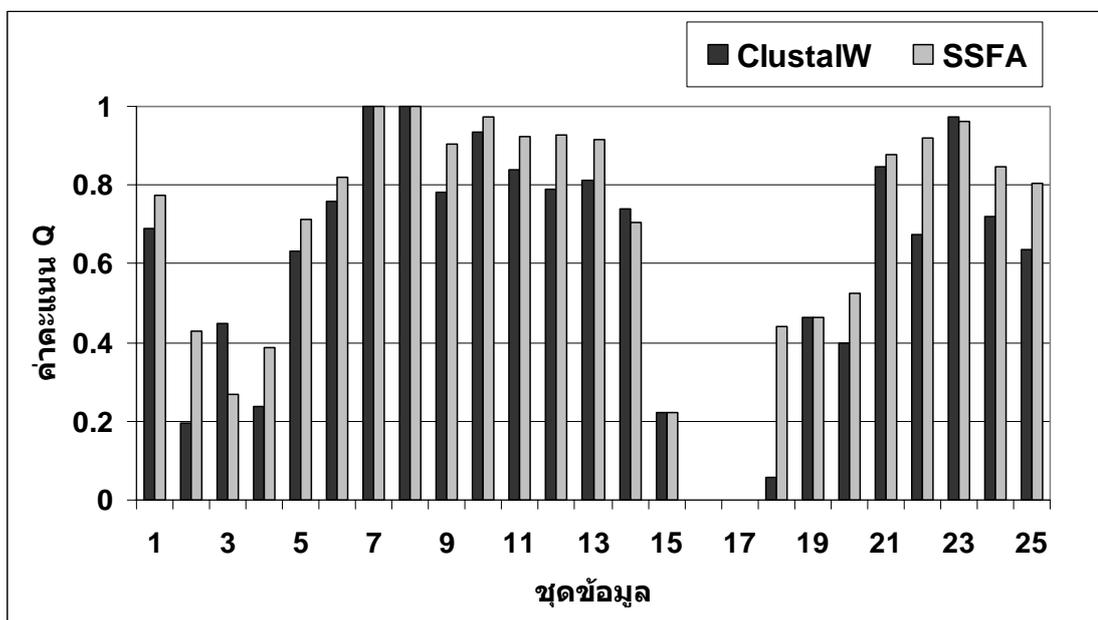
ภาพผนวกที่ 25 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 25



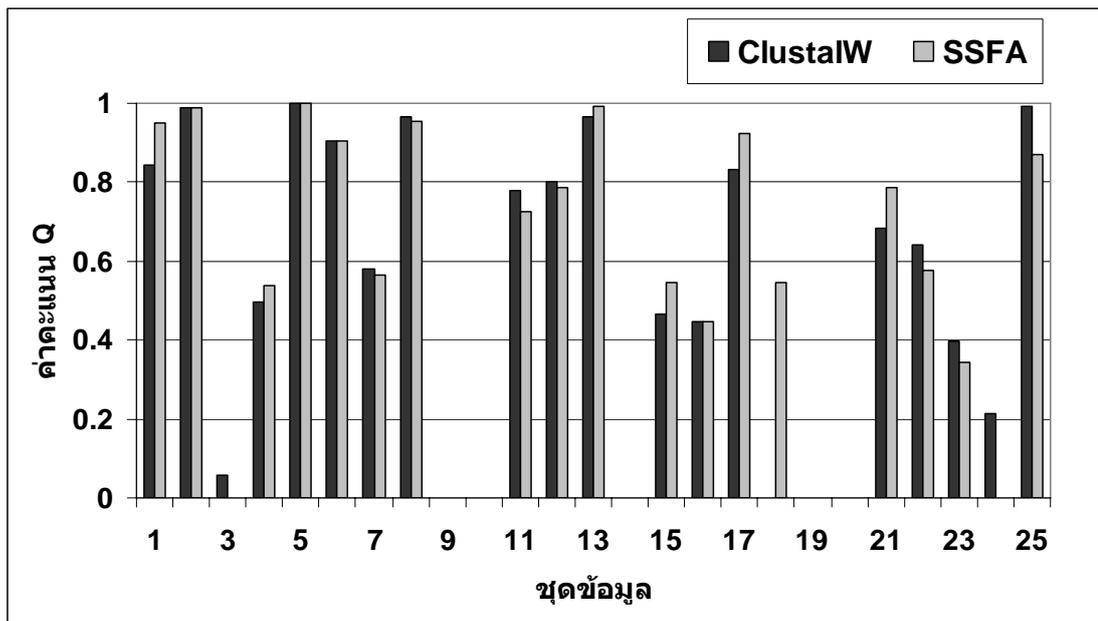
ภาพผนวกที่ 26 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 26



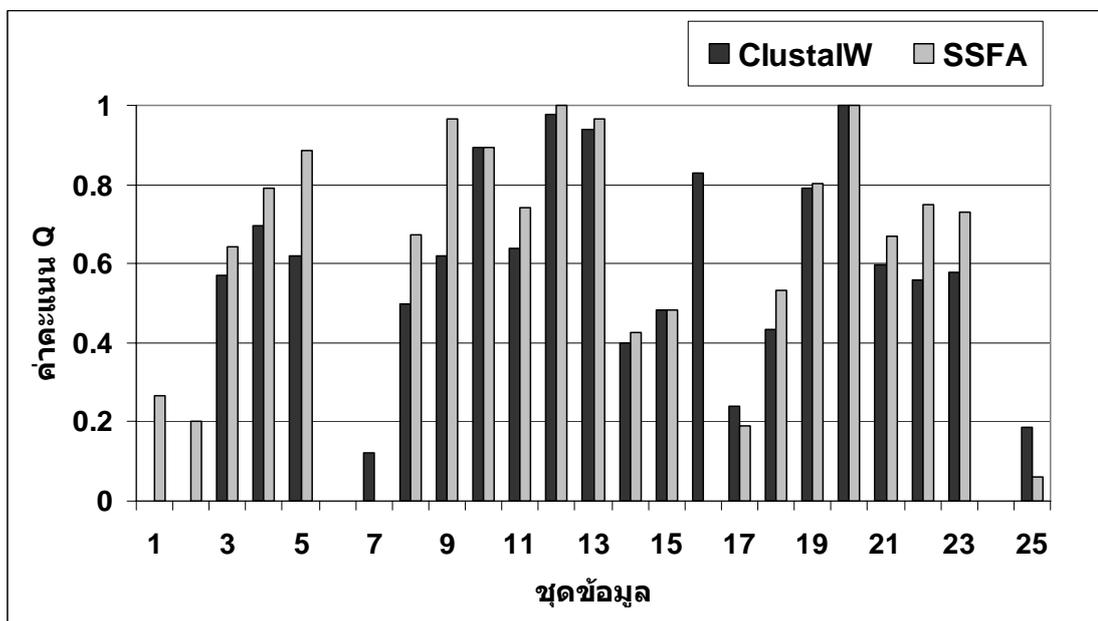
ภาพผนวกที่ 27 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 27



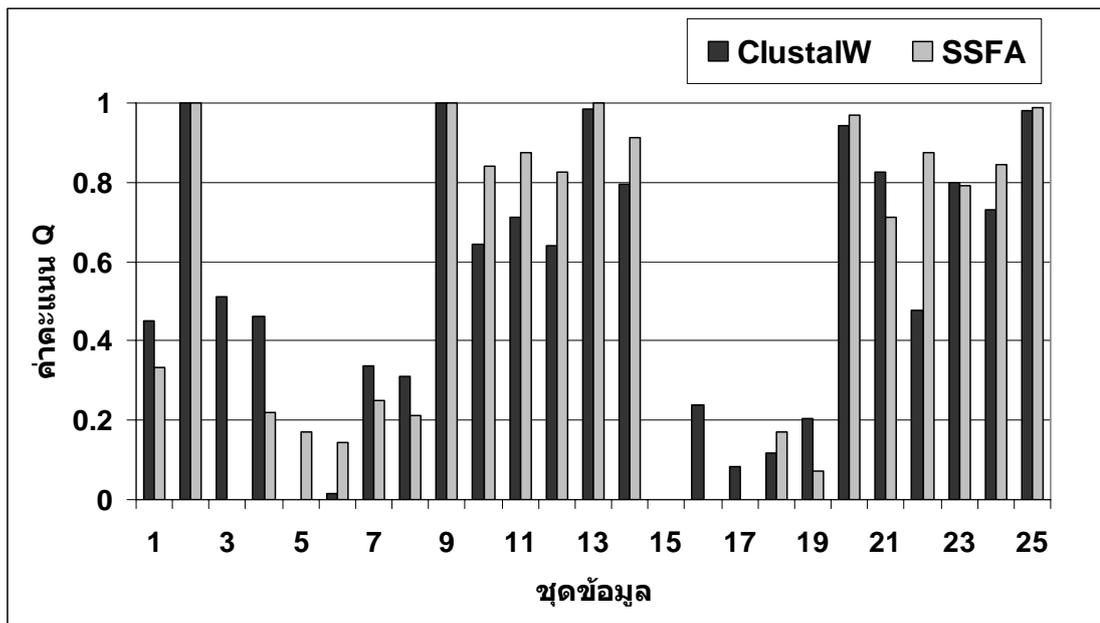
ภาพผนวกที่ 28 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 28



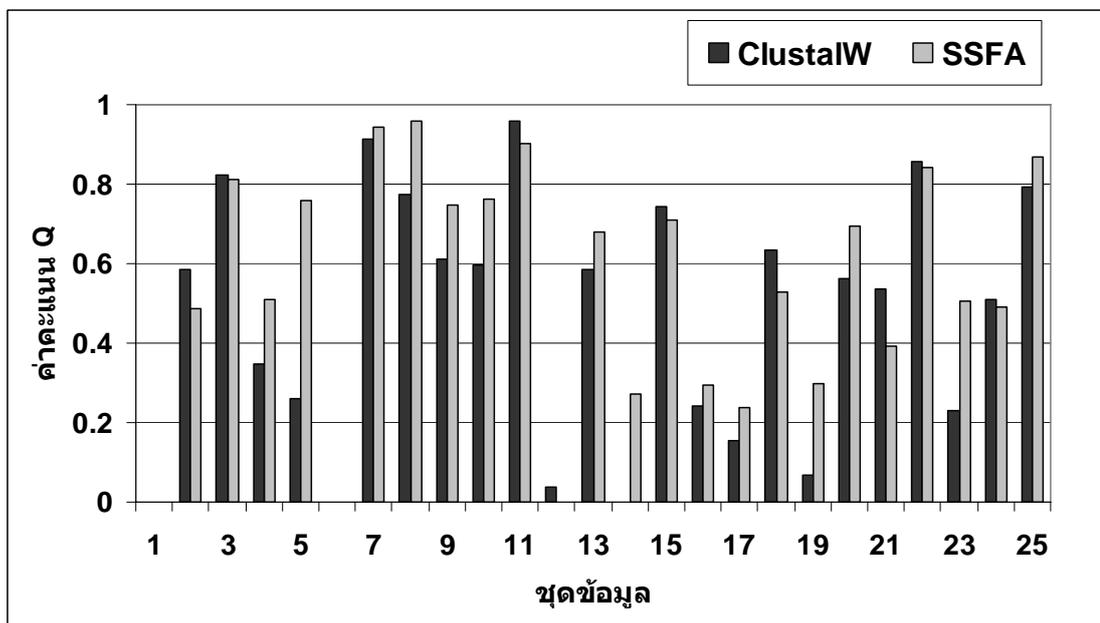
ภาพผนวกที่ 29 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 29



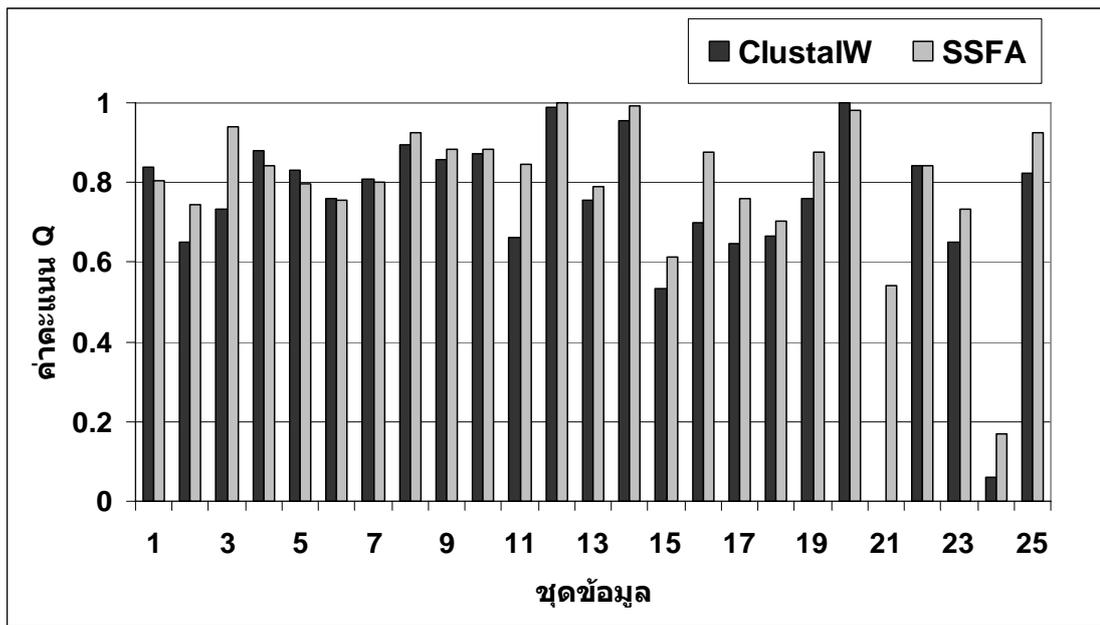
ภาพผนวกที่ 30 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 30



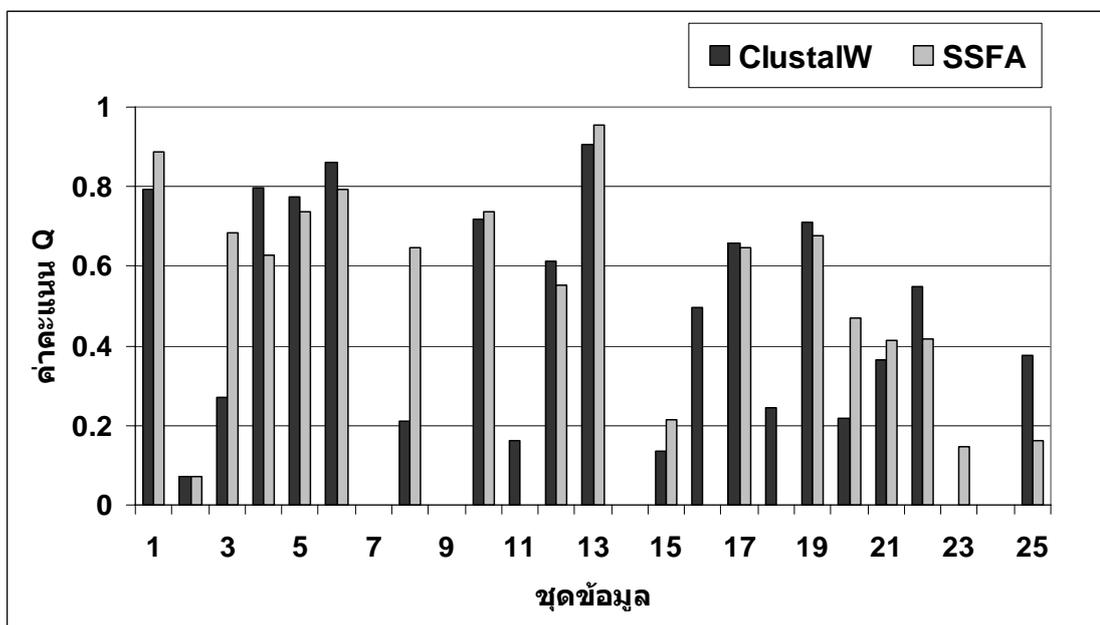
ภาพผนวกที่ 31 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 31



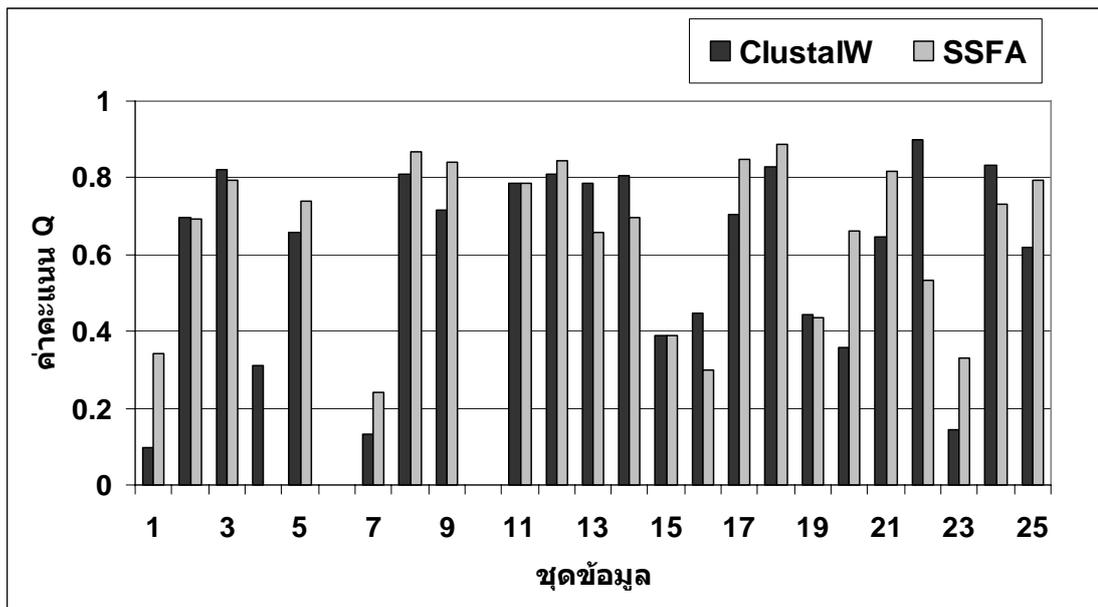
ภาพผนวกที่ 32 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 32



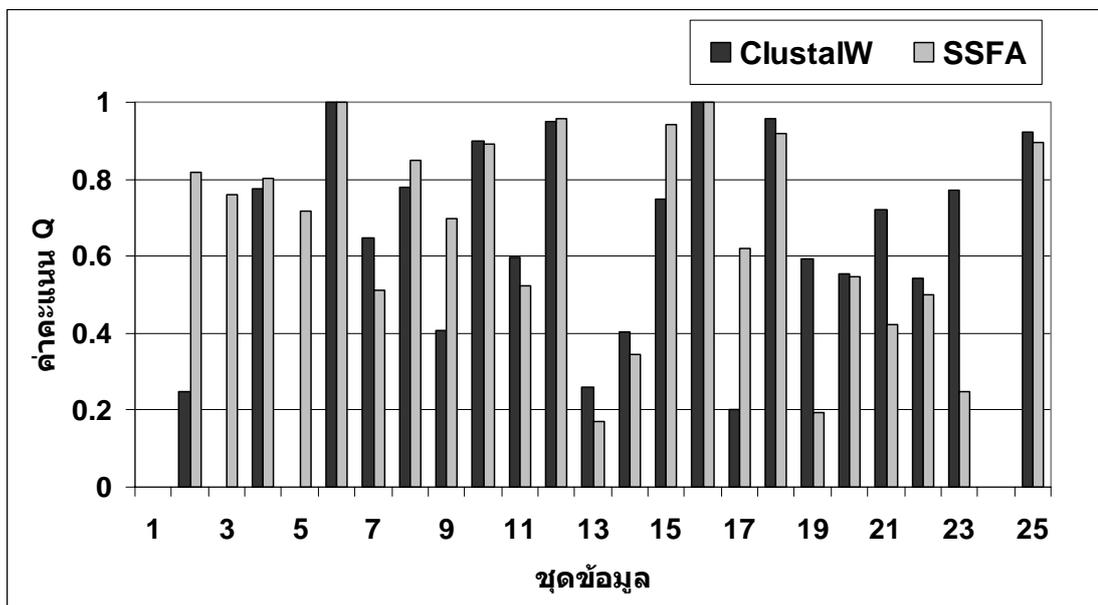
ภาพผนวกที่ 33 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของ โปรแกรม ClustalW และ โปรแกรม SSFA กับกลุ่มข้อมูลที่ 33



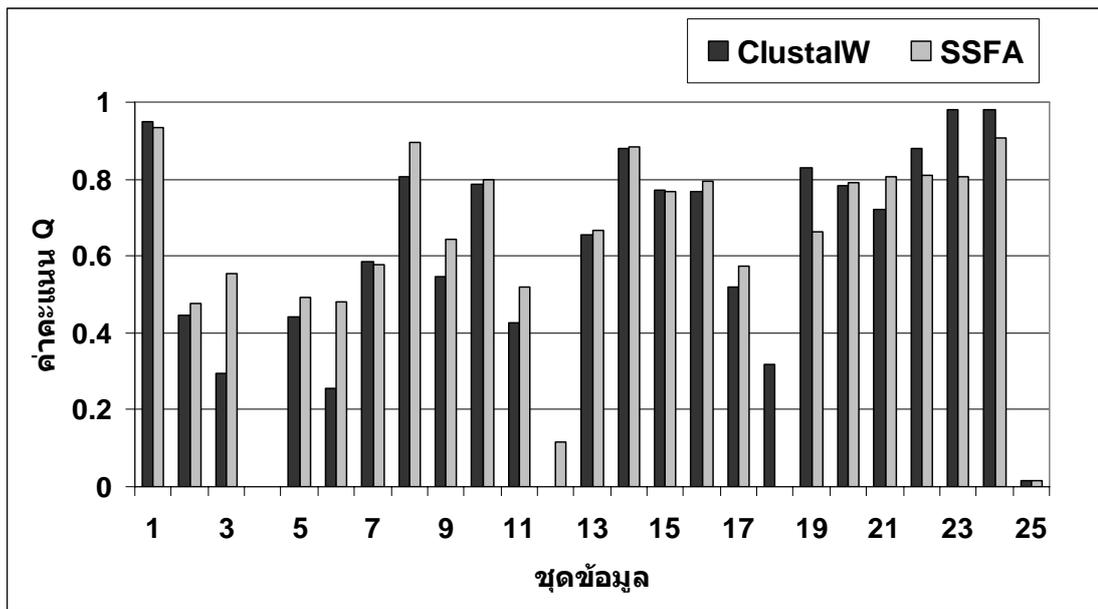
ภาพผนวกที่ 34 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของ โปรแกรม ClustalW และ โปรแกรม SSFA กับกลุ่มข้อมูลที่ 34



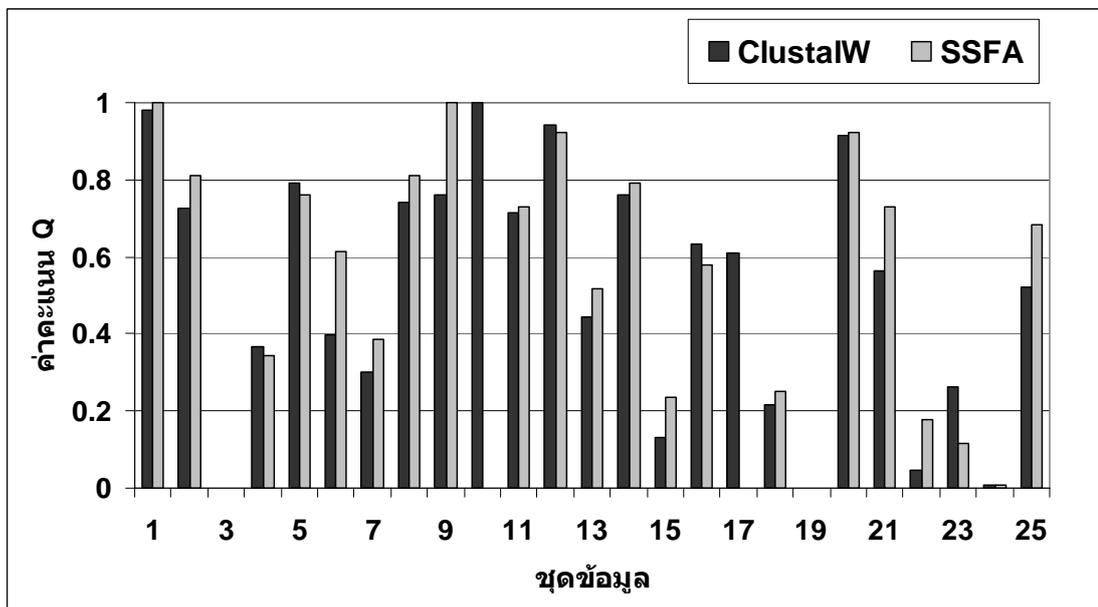
ภาพผนวกที่ 35 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 35



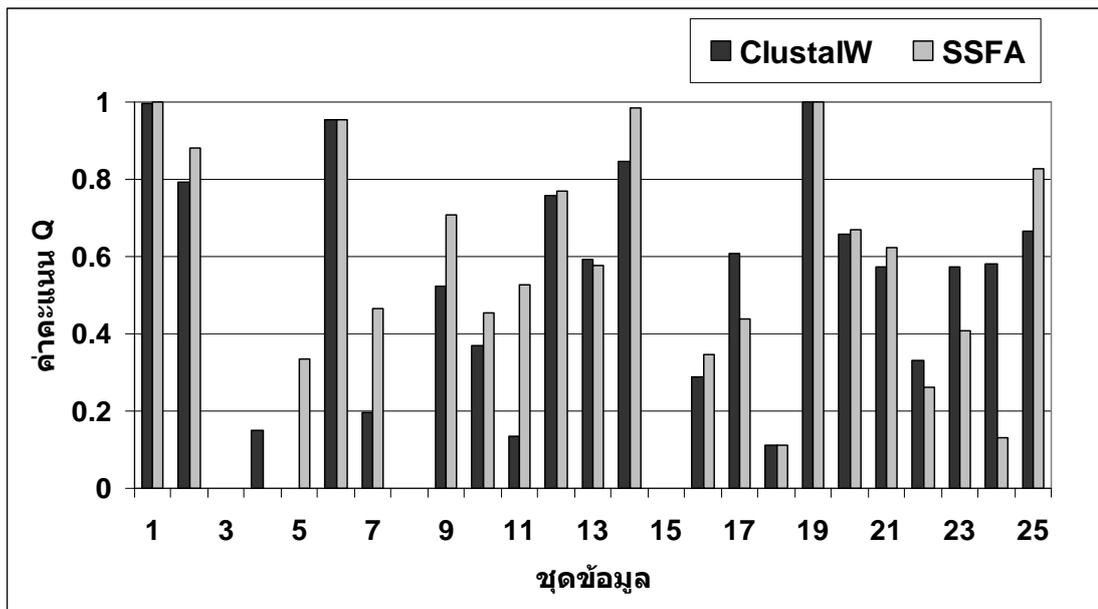
ภาพผนวกที่ 36 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 36



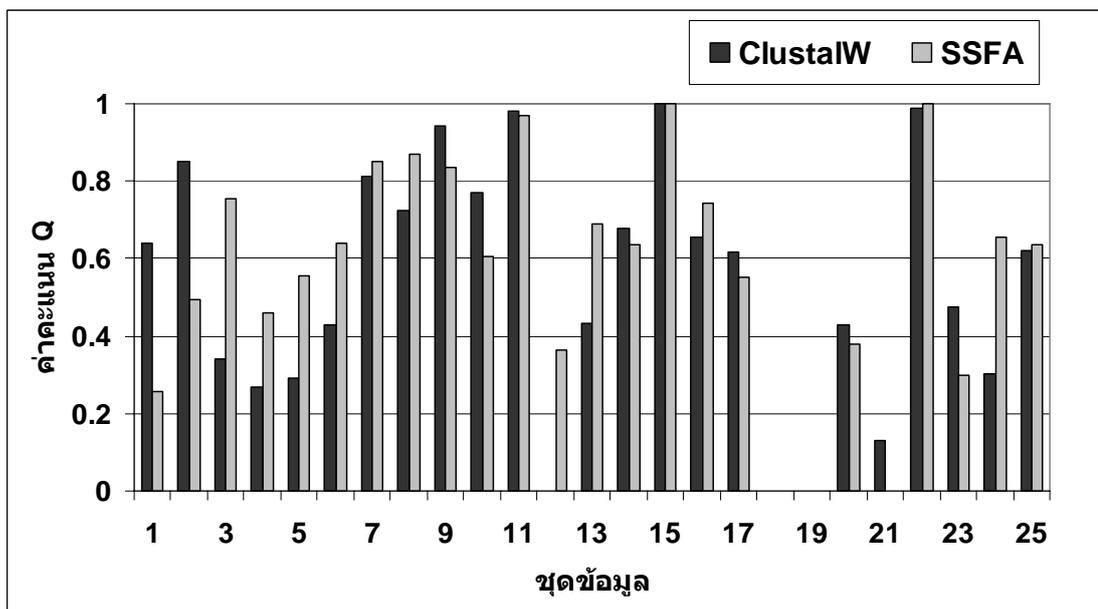
ภาพผนวกที่ 37 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 37



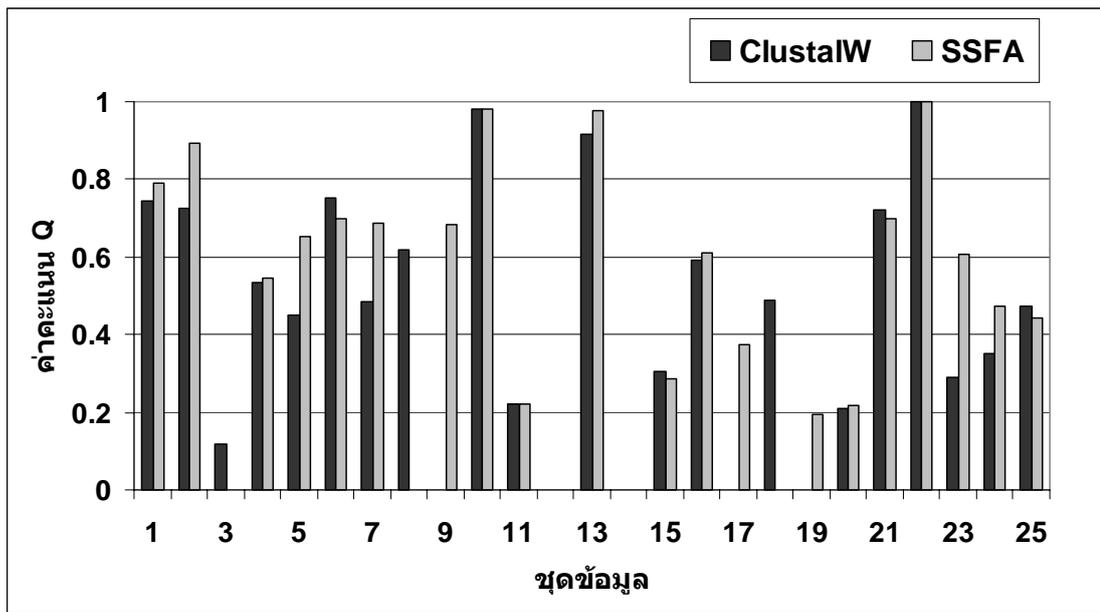
ภาพผนวกที่ 38 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 38



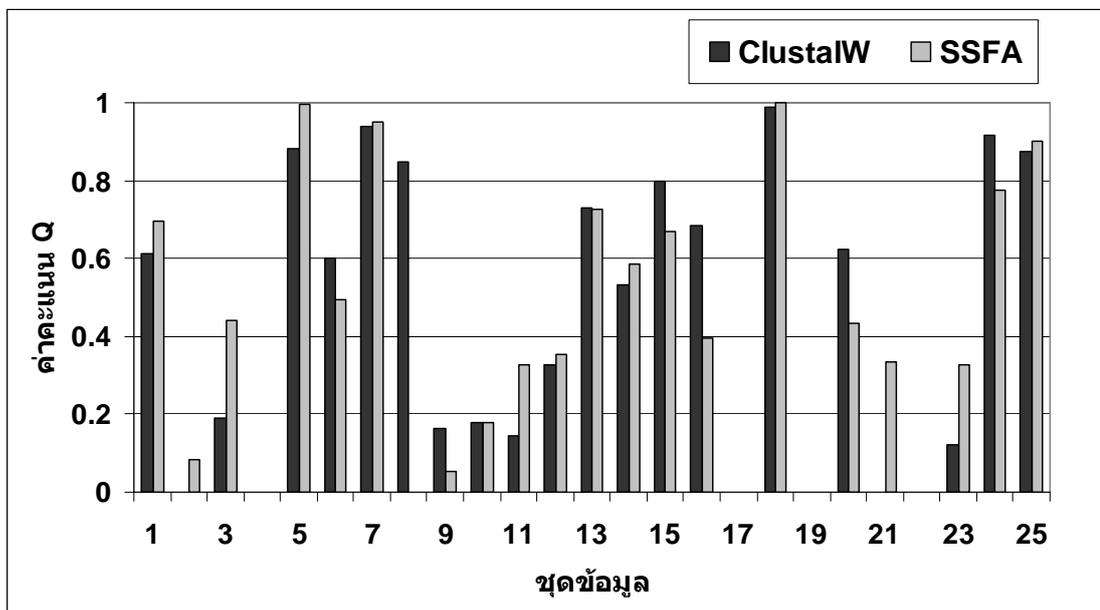
ภาพผนวกที่ 39 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 39



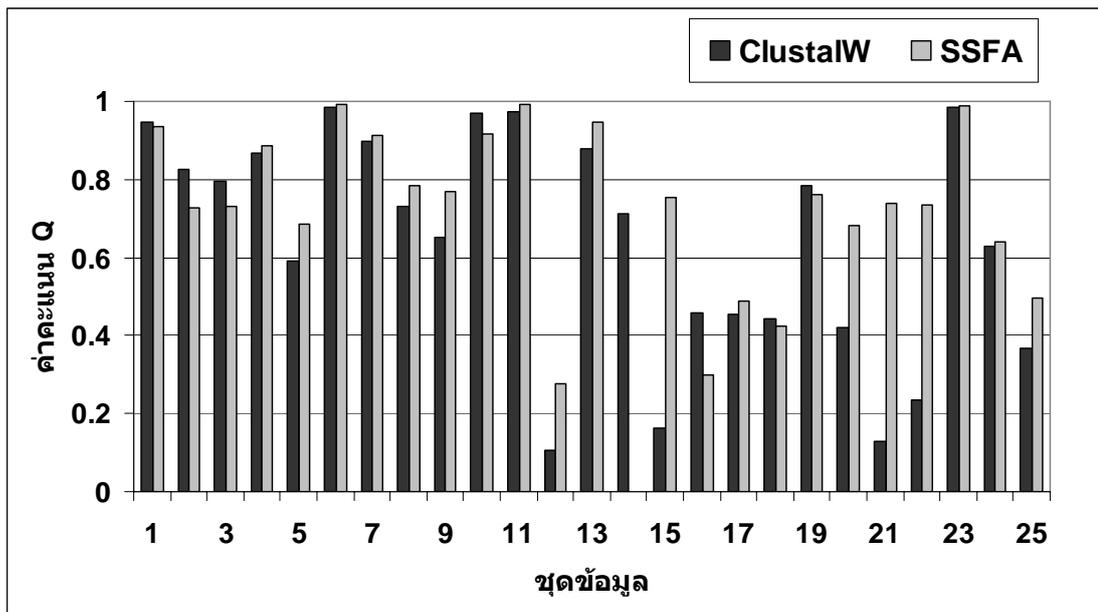
ภาพผนวกที่ 40 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 40



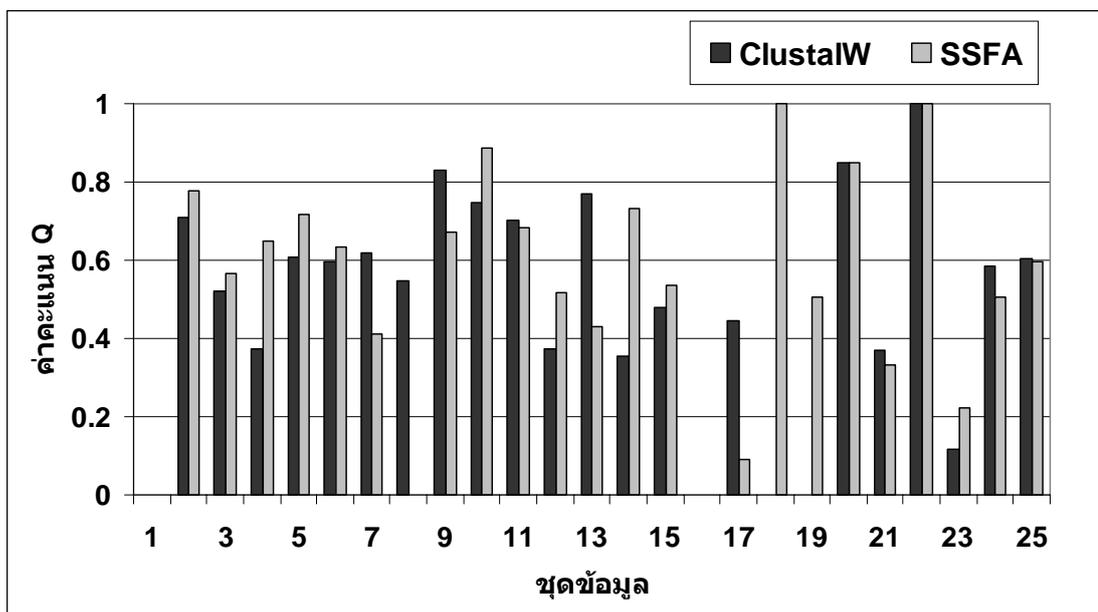
ภาพผนวกที่ 41 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 41



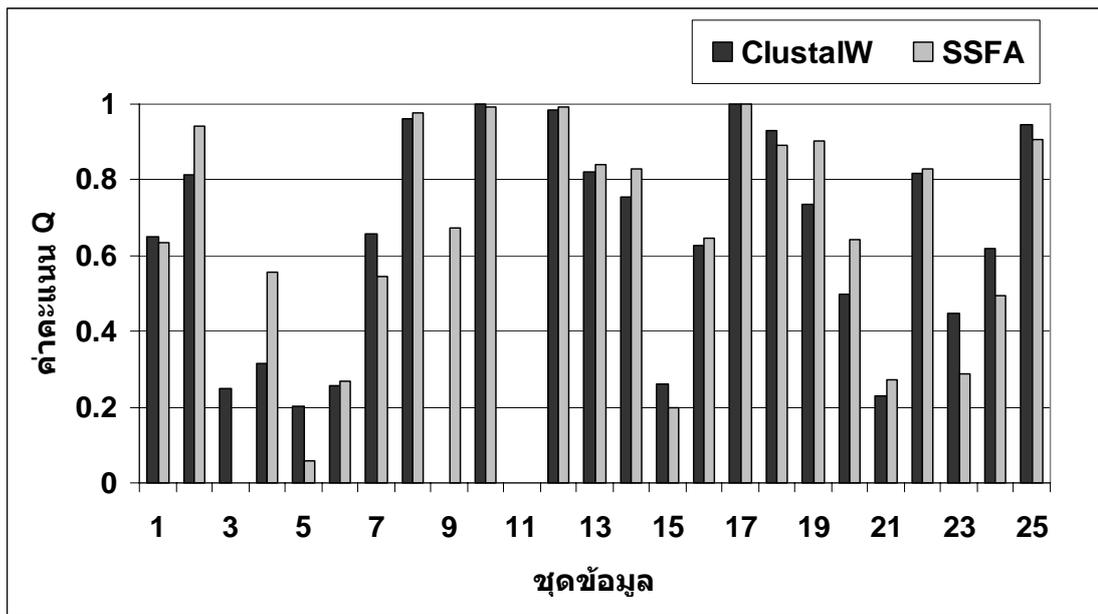
ภาพผนวกที่ 42 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 42



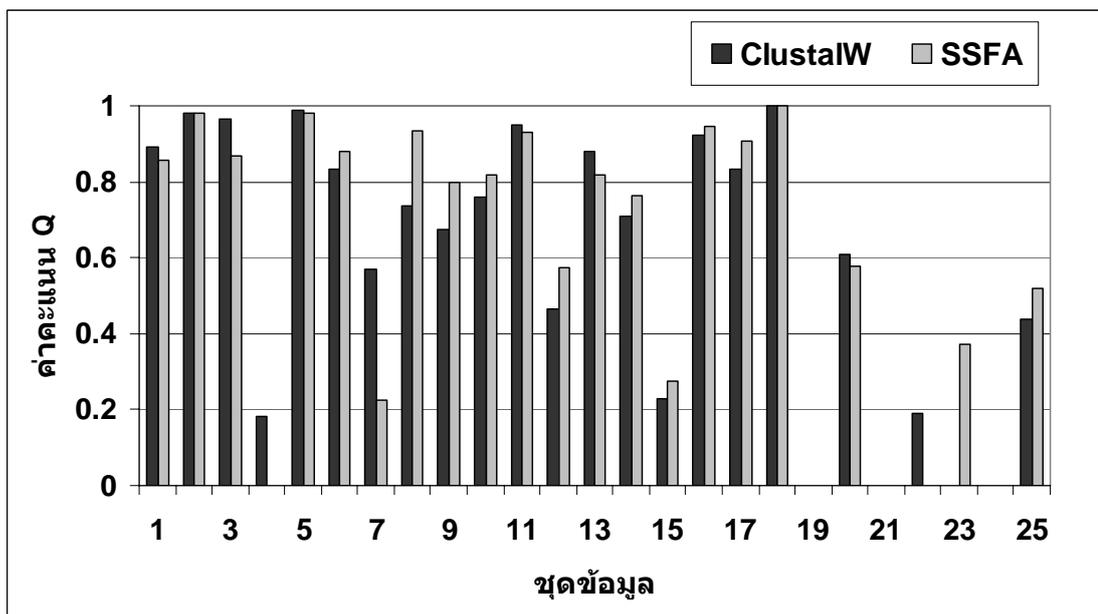
ภาพผนวกที่ 43 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 43



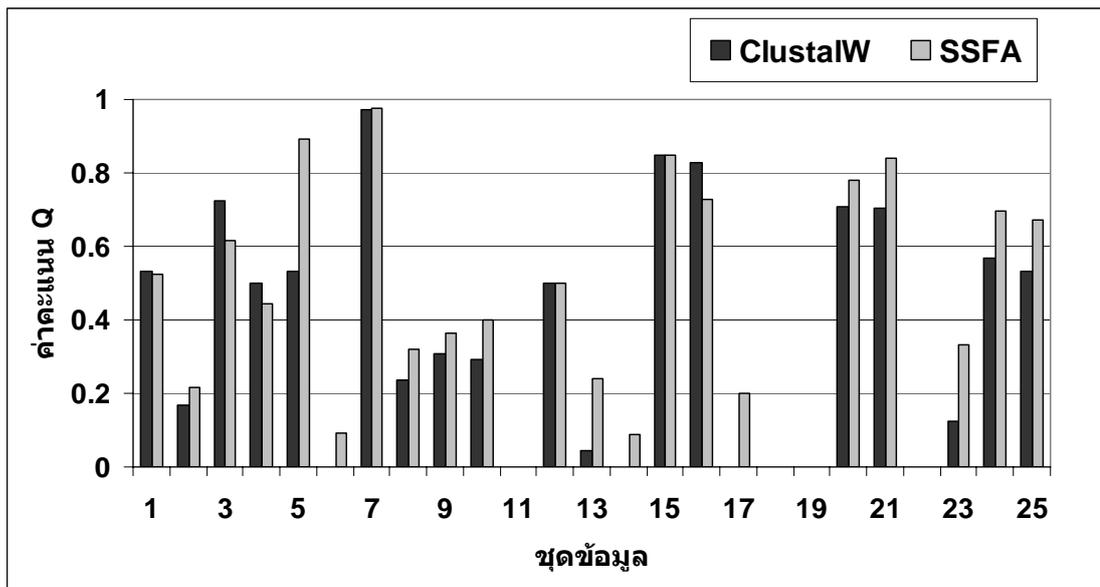
ภาพผนวกที่ 44 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 44



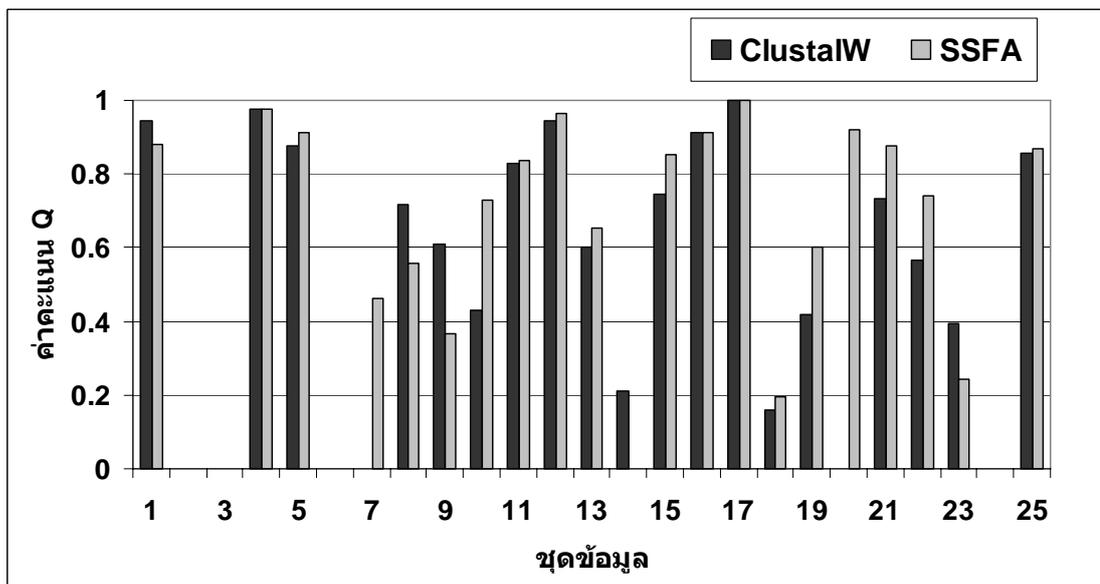
ภาพผนวกที่ 45 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 45



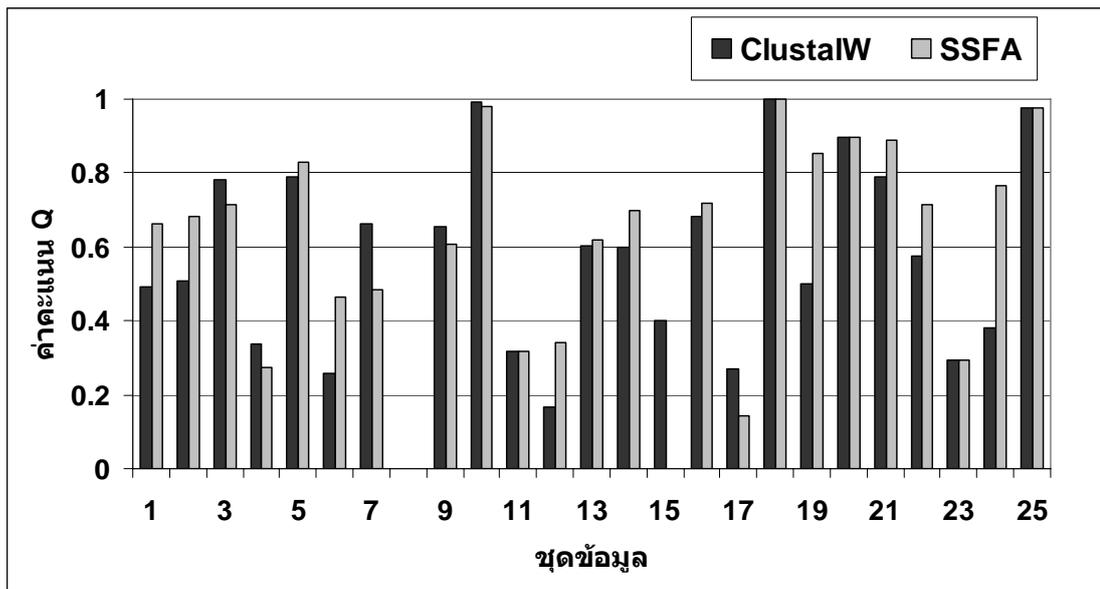
ภาพผนวกที่ 46 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 46



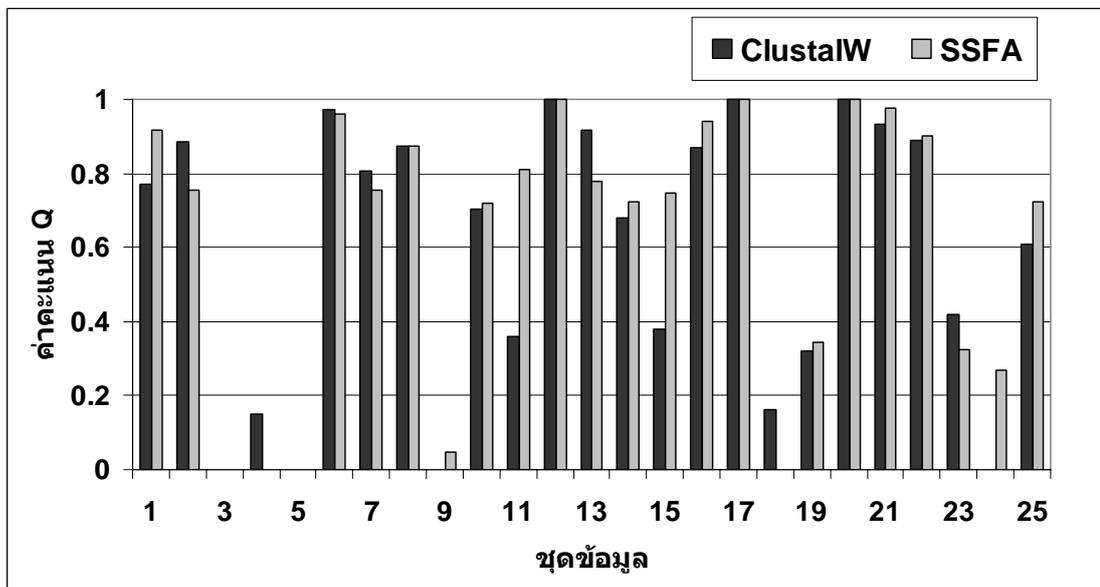
ภาพผนวกที่ 47 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 47



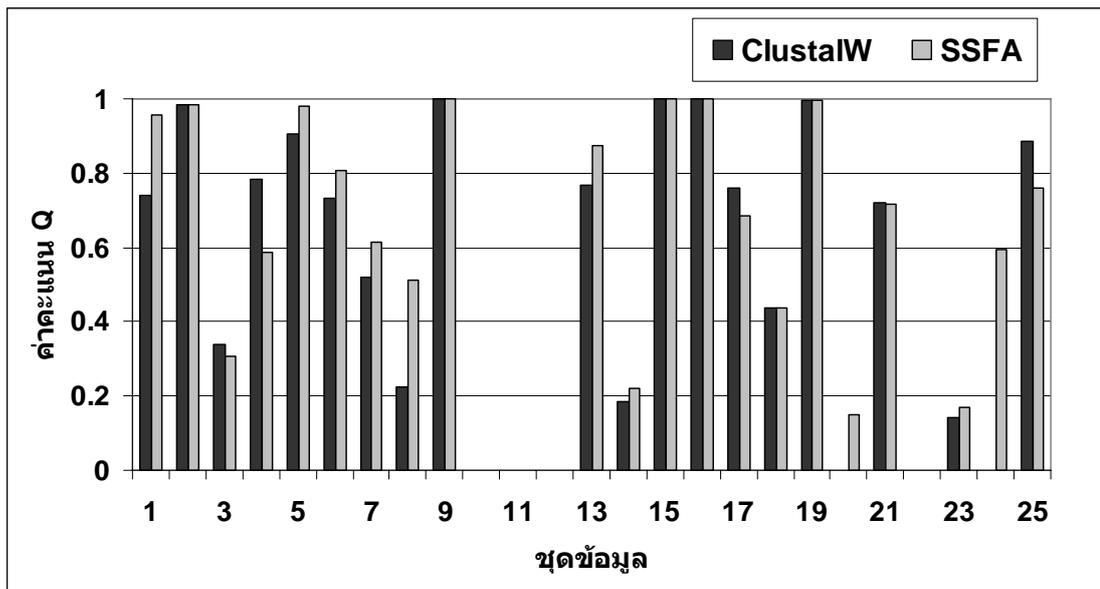
ภาพผนวกที่ 48 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 48



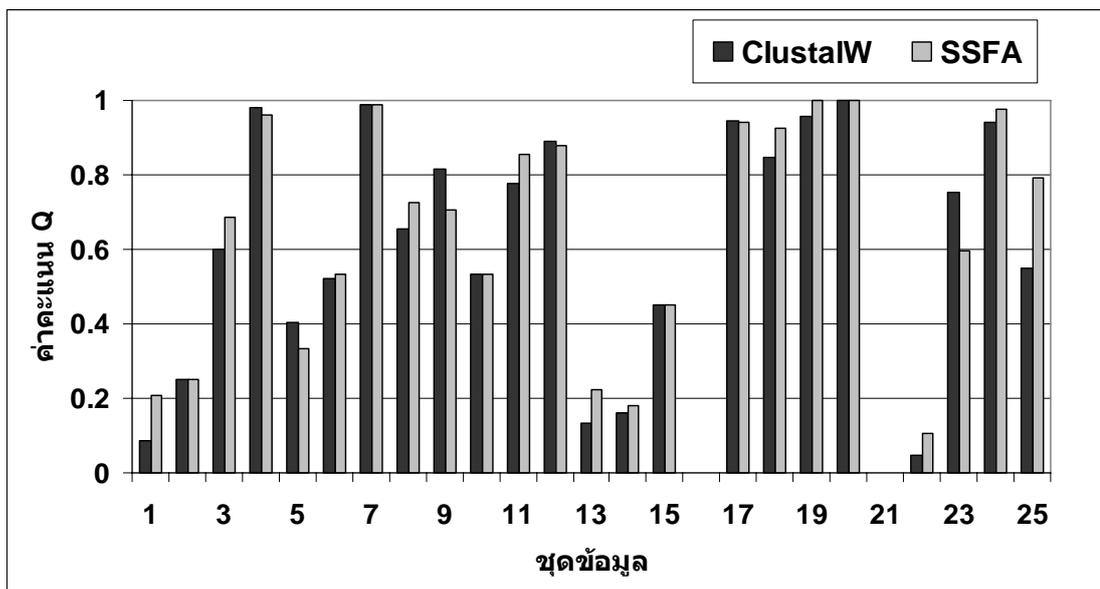
ภาพผนวกที่ 49 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 49



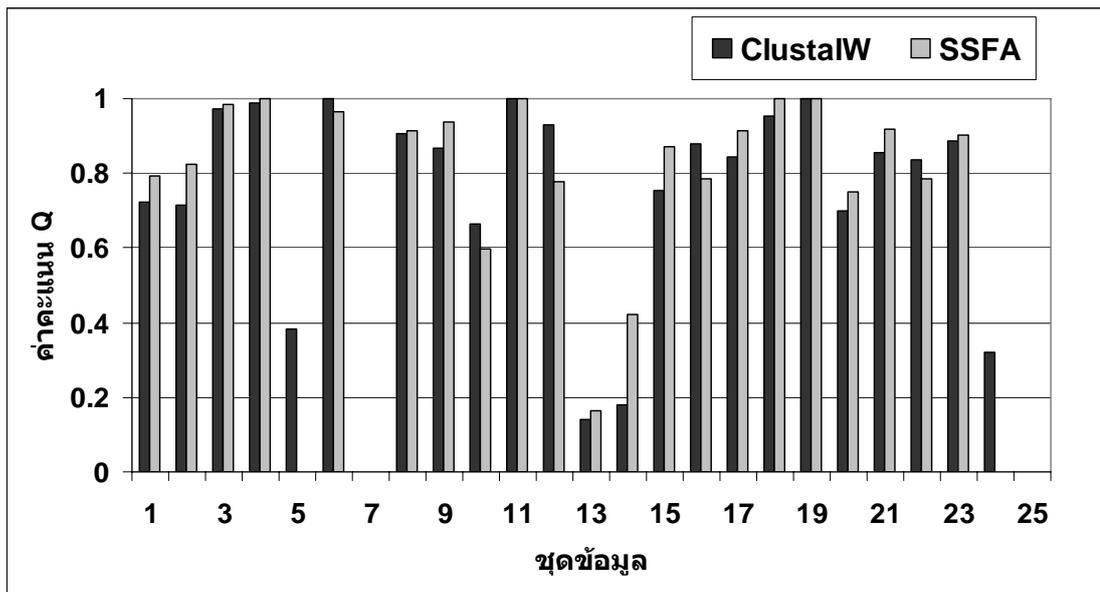
ภาพผนวกที่ 50 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 50



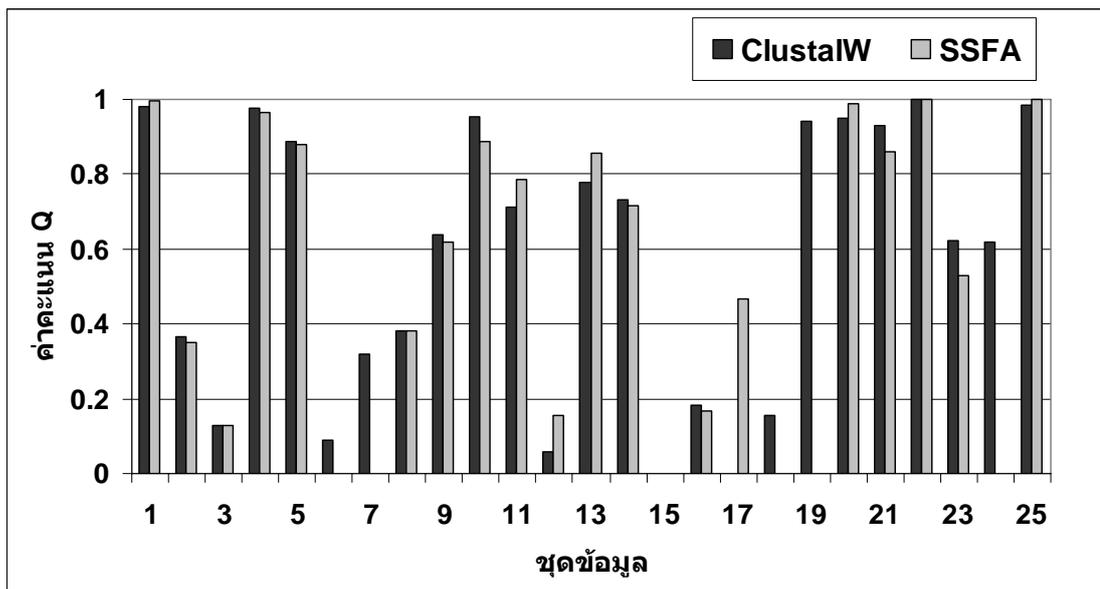
ภาพผนวกที่ 51 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 51



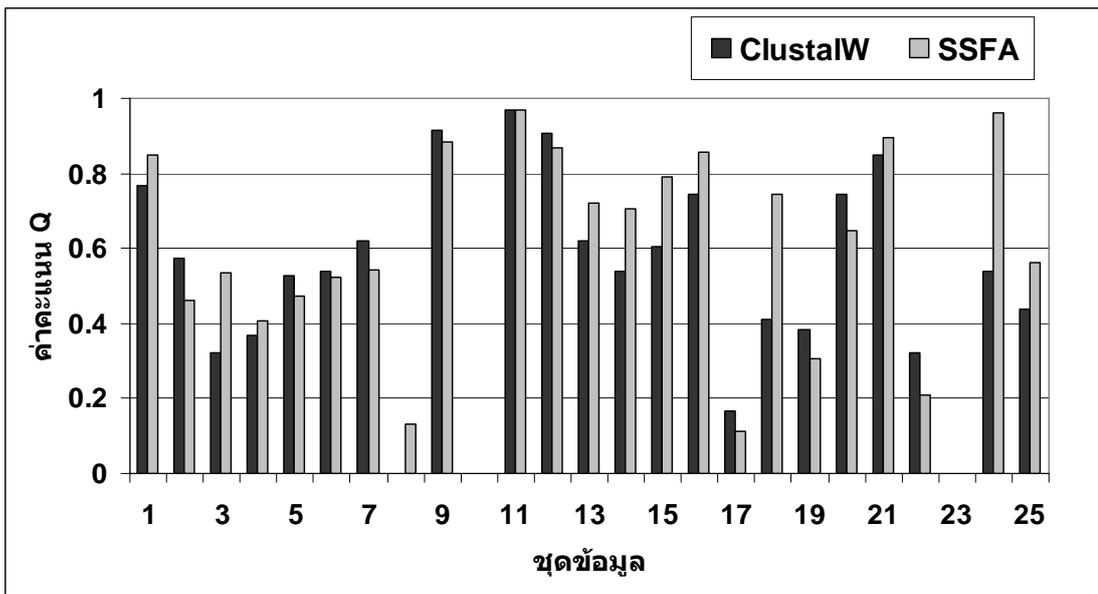
ภาพผนวกที่ 52 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 52



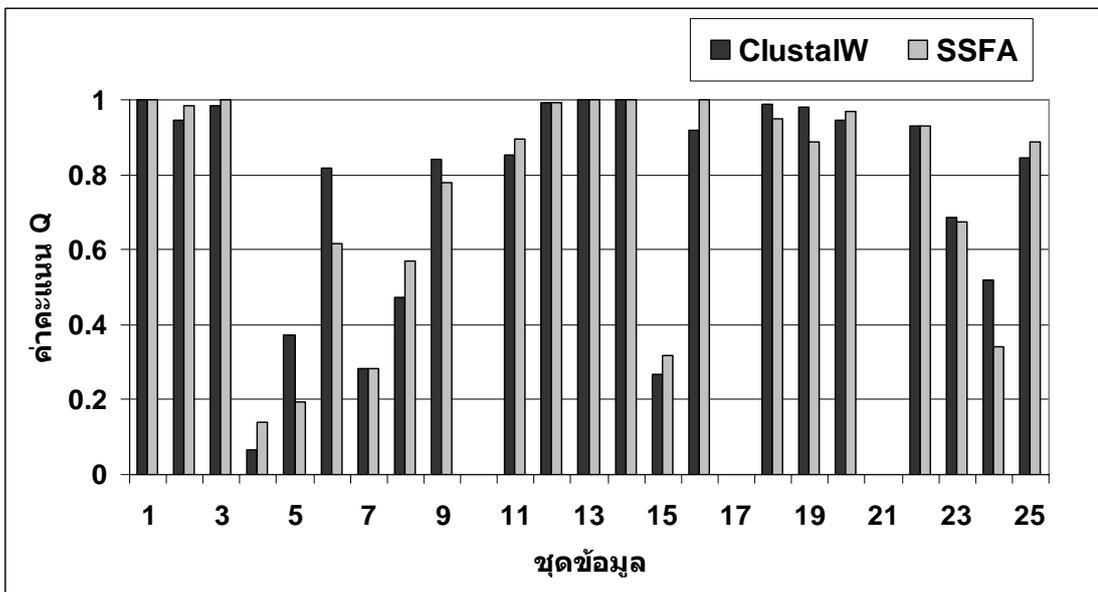
ภาพผนวกที่ 53 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 53



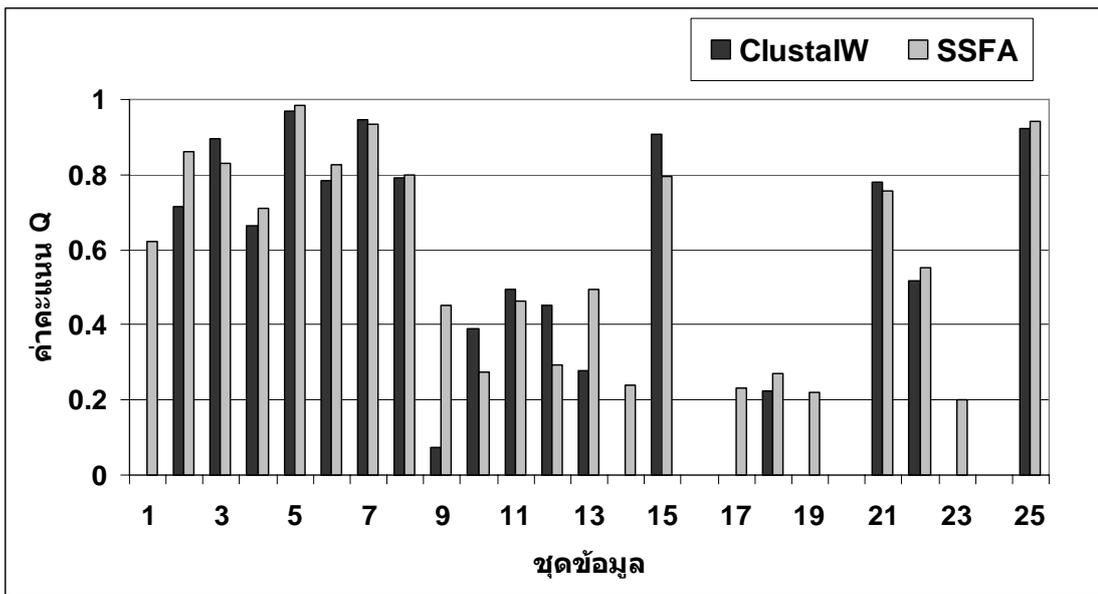
ภาพผนวกที่ 54 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 54



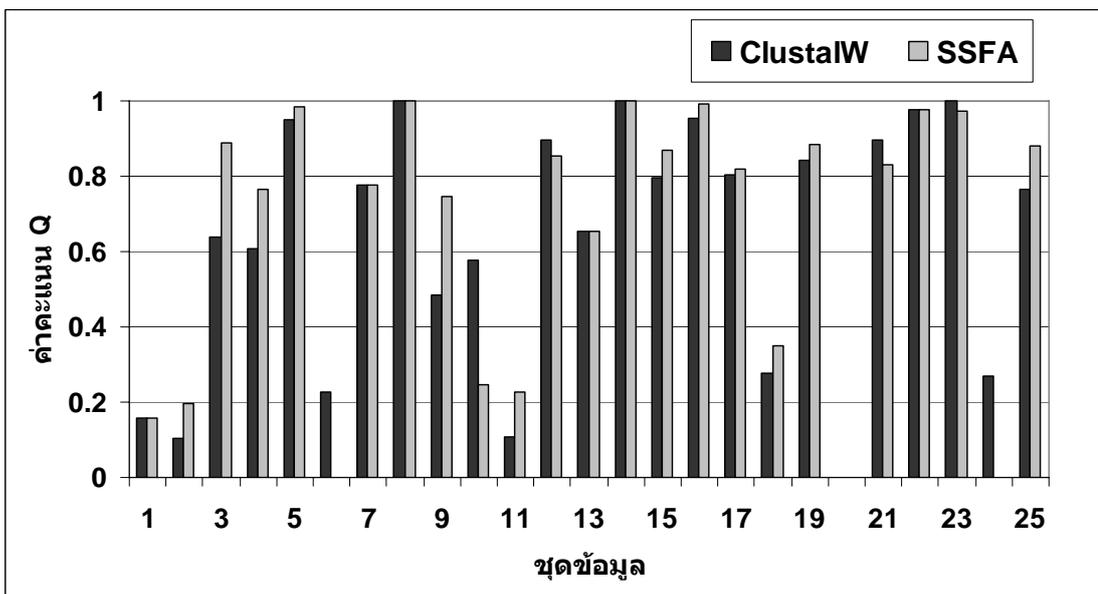
ภาพผนวกที่ 55 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของ โปรแกรม ClustalW และ โปรแกรม SSFA กับกลุ่มข้อมูลที่ 55



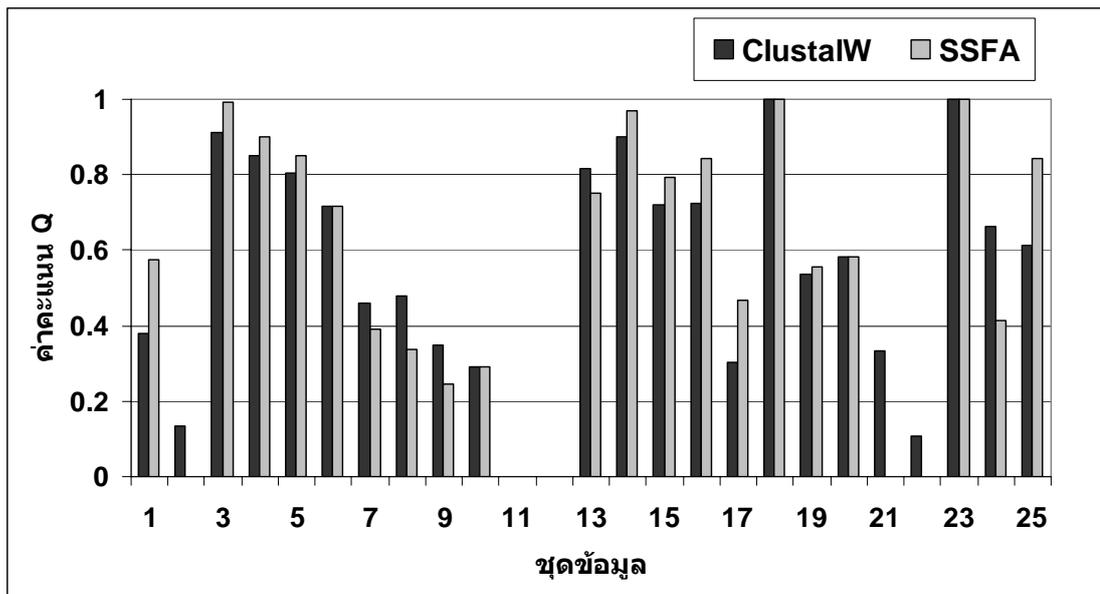
ภาพผนวกที่ 56 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของ โปรแกรม ClustalW และ โปรแกรม SSFA กับกลุ่มข้อมูลที่ 56



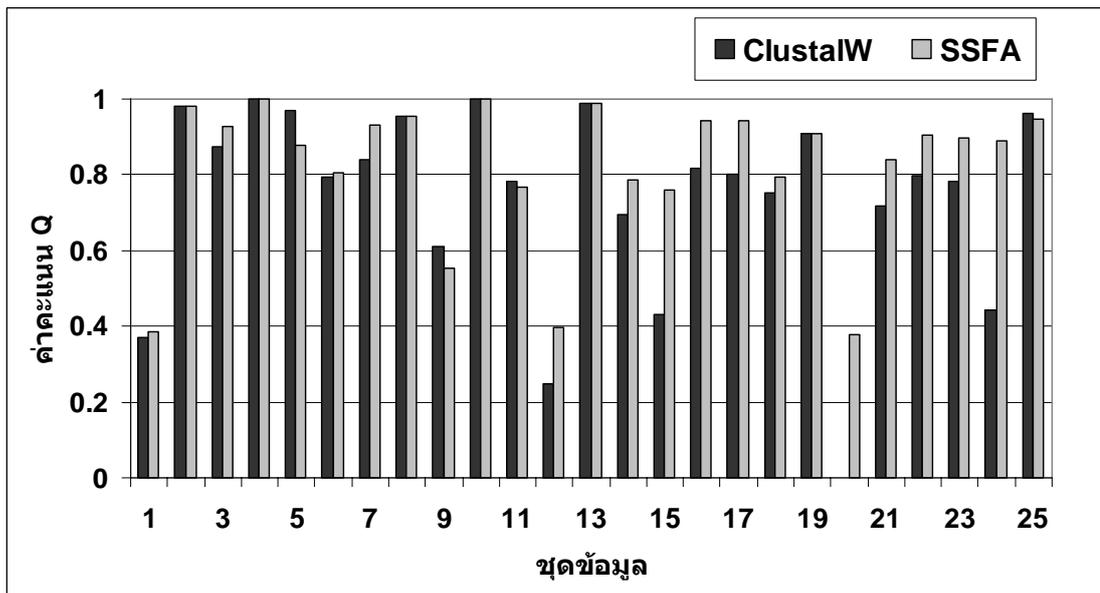
ภาพผนวกที่ 57 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 57



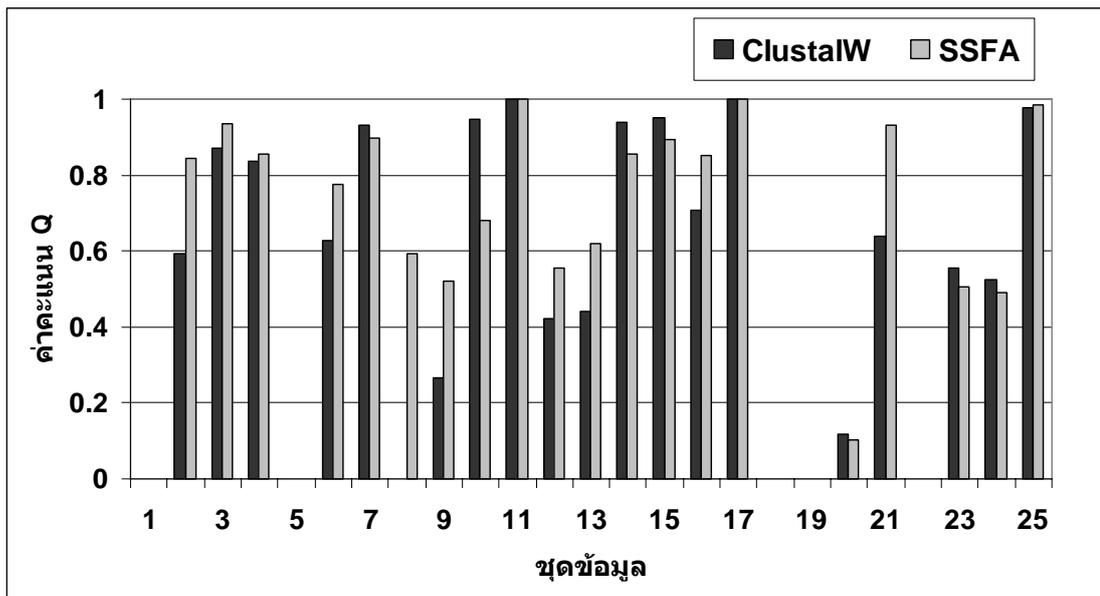
ภาพผนวกที่ 58 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 58



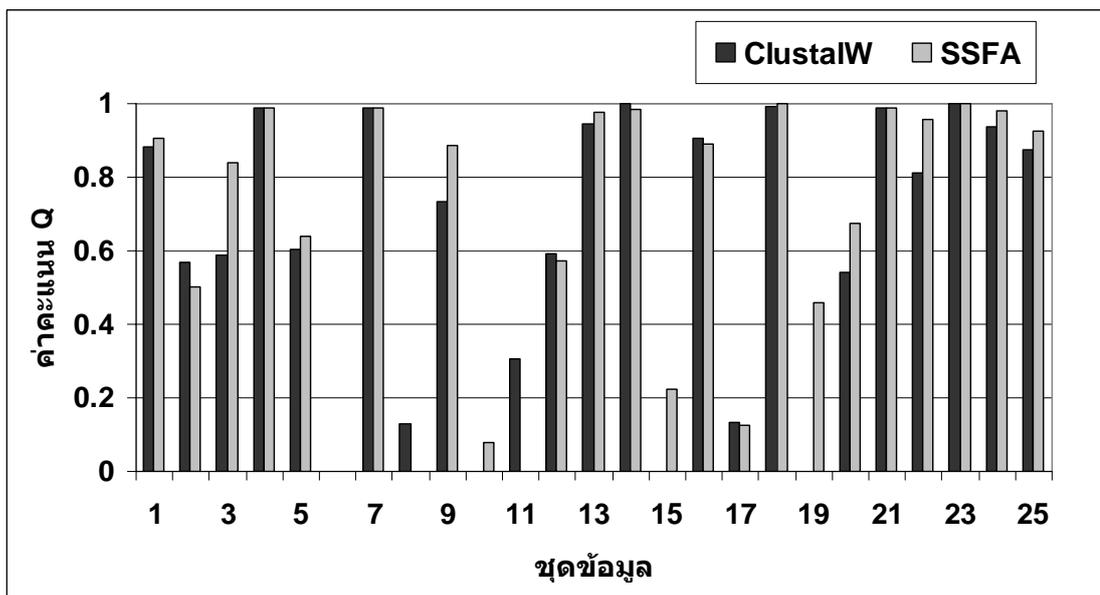
ภาพผนวกที่ 59 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 59



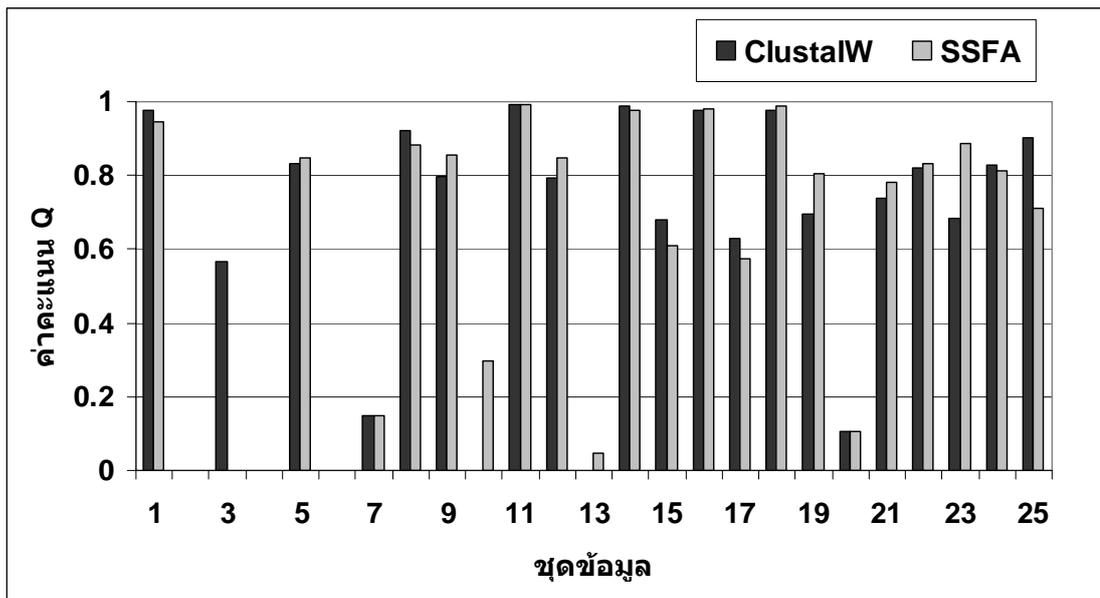
ภาพผนวกที่ 60 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 60



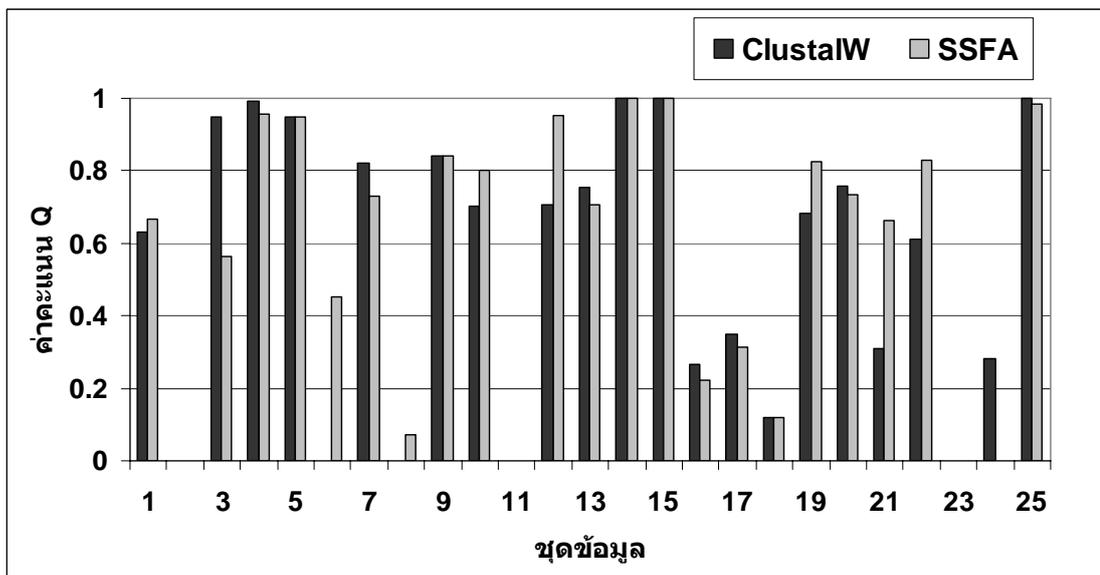
ภาพผนวกที่ 61 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 61



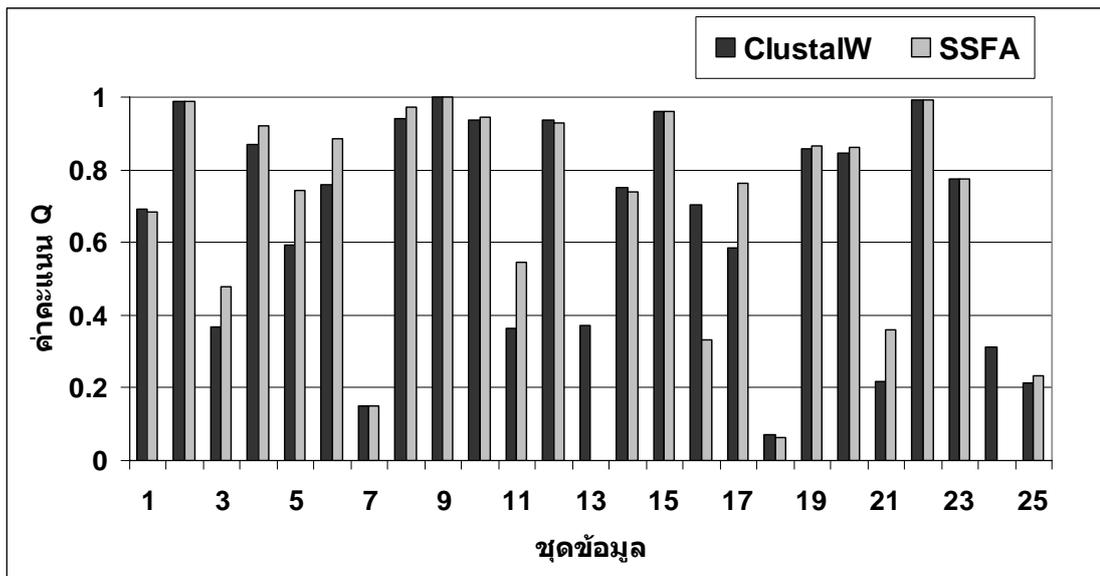
ภาพผนวกที่ 62 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 62



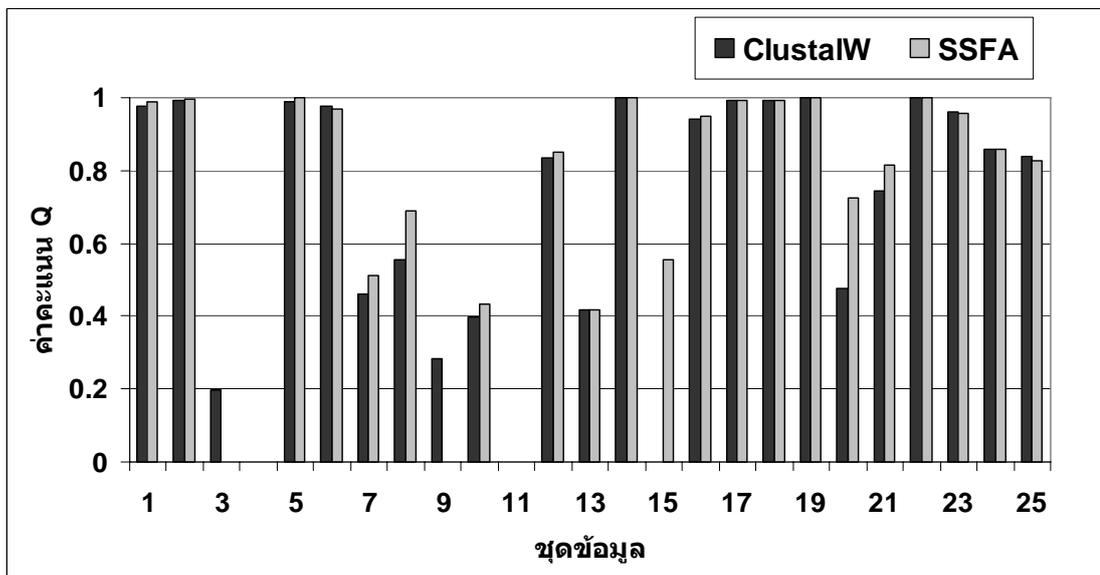
ภาพผนวกที่ 63 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 63



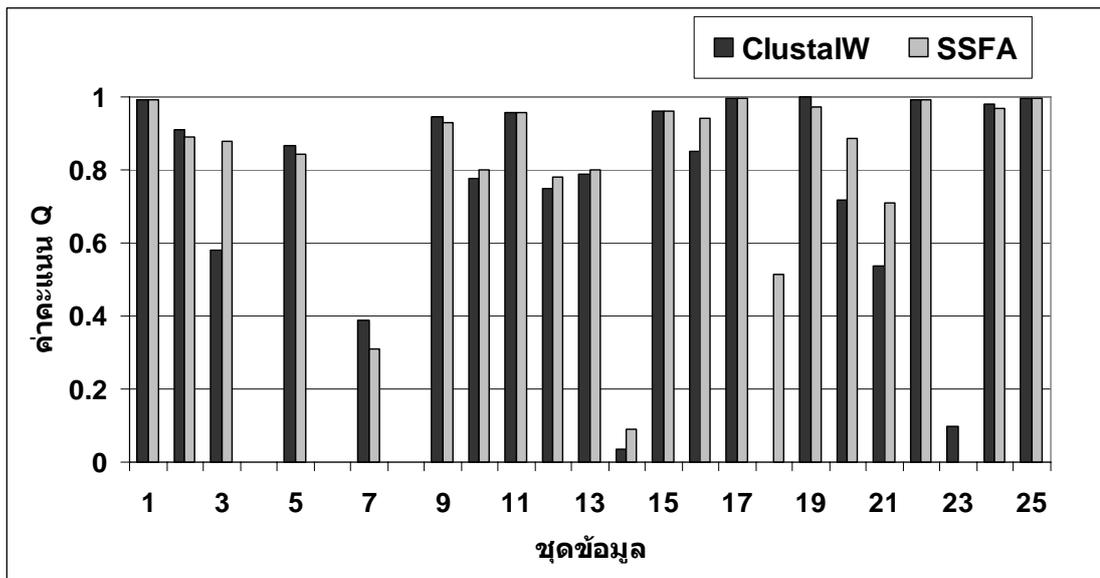
ภาพผนวกที่ 64 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 64



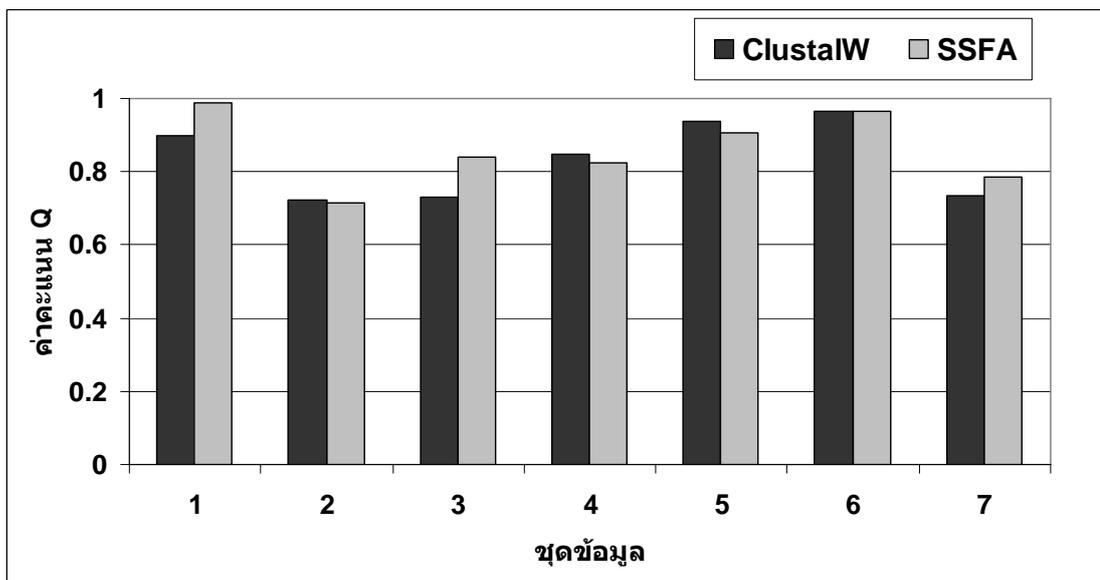
ภาพผนวกที่ 65 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 65



ภาพผนวกที่ 66 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 66



ภาพผนวกที่ 67 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 67



ภาพผนวกที่ 68 ค่าคะแนน Q ในการเทียบเรียงฐานข้อมูล PREFAB ของโปรแกรม ClustalW และโปรแกรม SSFA กับกลุ่มข้อมูลที่ 68

ประวัติการศึกษา และการทำงาน

ชื่อ –นามสกุล	นายคมสัน จันมา
วัน เดือน ปี ที่เกิด	14 พฤศจิกายน 2523
สถานที่เกิด	จังหวัดกรุงเทพมหานคร
ประวัติการศึกษา	ปริญญาวิศวกรรมศาสตรบัณฑิต (วิศวกรรมคอมพิวเตอร์)
ตำแหน่งหน้าที่การงานปัจจุบัน	วิศวกรระบบ
สถานที่ทำงานปัจจุบัน	บริษัท ที.เอ็น. อินฟอร์เมชั่น ซิสเต็มส์ จำกัด
ผลงานดีเด่นและรางวัลทางวิชาการ	-
ทุนการศึกษาที่ได้รับ	-