

ขั้นตอนวิธีการตัดคำ  
สำหรับการตรวจสอบตัวสะกดภาษาไทย

**A word segmentation algorithm  
for Thai spell checking**

ทิวา เจริญสวัสดิพงศ์  
ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์  
จุฬาลงกรณ์มหาวิทยาลัย

# ปัญหา

- ปัญหาของการตรวจสอบตัวสะกดภาษาไทย คือ เขียนคำติดกัน
- ไม่ทราบขอบเขตของคำ
- ต้องนำข้อความมาผ่านการตัดคำก่อน
- นำคำที่ได้ไปเทียบกับคำในพจนานุกรม

# เสนอ

- อัลกอริทึมการตัดคำเพื่อการตรวจสอบตัวสะกดภาษาไทย
- ประสิทธิภาพการตรวจสอบตัวสะกดดีกว่าที่มีอยู่ในปัจจุบัน

# โปรแกรมประมวลผลคำทั่วไป

- ไมโครซอฟท์เวิร์ด 2002, ปลาตาวออฟฟิศ 2.0
- จีดเส้นใต้เพื่อระบุการสะกดผิด ไม่ตรงกับจุดที่ผิดจริง ๆ ซึ่งมีผลทำให้ไม่สามารถแก้ไขคำนั้นได้ถูกต้อง

คำที่สะกดผิด	สะกดผิดจาก	ตำแหน่งที่ผิด	บริเวณที่จีดเส้นใต้ระบุการสะกดผิด	
			ไมโครซอฟท์เวิร์ด 2002	ปลาตาวออฟฟิศ 2.0
อิเล็กทรอนิกส์	อิเล็กทรอนิกส์	ต (พิมพ์เกิน)	อิเล็กทรอนิกส์	อิเล็กทรอนิกส์
วิชาการ	วิชาการ	ก (พิมพ์สลับ)	วิชาการ	วิชาการ
ไผ่ยากรณ์	ไวยากรณ์	ฝ (พิมพ์ผิด)	กรณ์	กรณ์
อุปกรณ์	อุปกรณ์	ณ (พิมพ์ตก)	อุปก	ไม่มี

อิเล็กทรอนิกส์

วิชาการ

ไผ่ยากรณ์

อุปกรณ์

## โปรแกรมประมวลผลคำทั่วไป (2)

- หากคำผิดสามารถแตกย่อยเป็นคำที่ถูกต้อง จะไม่ระบุว่าเป็นคำผิด

คำที่สะกดผิด	พิมพ์ผิดจาก	คำย่อย
รัฐมนตรี	รัฐมนตรี	รัฐ นม ตรี
รัฐปะทาร์	รัฐปะทาร์	รัฐ ปะ ทาร์
สามทาร์ล	สามทาร์ล	สาม ทา รล

# โปรแกรมตัดคำทั่วไป

- โปรแกรมที่นำมาทดสอบ ได้แก่ Swath, ICU
- ตัดคำผิดพลาดเมื่อได้รับข้อความที่มีคำสะกดผิด

คำที่สะกดผิด	สะกดผิดจาก	ผลการตัดคำจาก โปรแกรมตัดคำ	
		Swath	ICU
เทคร โน โสยี่	เทค โน โสยี่	เท คร  โน  โส ยี่	เท คร โน โส ยี่
เจตนรรมณั์	เจตนารมณั์	เจ ต น ร ม ณั์	เจ ต น ร ม ณั์

- คำยาวสะกดผิด แต่สามารถตัดเป็นคำถูกต้องได้ -> คำถูกย่อย ๆ

คำที่สะกดผิด	สะกดผิดจาก	ผลการตัดคำจาก โปรแกรมตัดคำ	
		Swath	ICU
รัฐนมนตรี	รัฐมนตรี	รัฐ น ม ต ร ี	รัฐ น ม ต ร ี
กระทรงวง	กระทรง	ก ระ ท ร ง ว ง	ก ระ ท ร ง ว ง

# การสะกดผิด

- คำที่พิมพ์ผิด เกิดขึ้นจากความผิดพลาด 4 ชนิด [Damerou, 1964 ]
  - พิมพ์เกิน เช่น วิทยาลัย พิมพ์เป็น วิทยาลัย
  - พิมพ์ตก เช่น อาหาร พิมพ์เป็น อหาร
  - พิมพ์ผิด เช่น นักเรียน พิมพ์เป็น นักรเือน
  - พิมพ์สลับ เช่น รัฐบาล พิมพ์เป็น รัฐบาล

# พฤติกรรมกรรมการพิมพ์ผิดพลาดของมนุษย์

- ผู้ที่พิมพ์ชำนาญมักผิดพลาดจากการพิมพ์เกินใน 2 แป้นที่อยู่ติดกัน
- ผู้ที่พิมพ์ไม่ชำนาญมักผิดพลาดจากการพิมพ์ผิด และการพิมพ์ผิดนั้นมักเกิดจากแป้นที่อยู่ติดกัน [Grudin, 1983]
- มีเพียงส่วนน้อยมากของความผิดพลาดที่เกิดจากอักขระตัวแรกของคำ [Kukich, 1992]

# วิธีการตัดค่าแบบเหมือนมากที่สุด (Maximal Matching) [วิรัช, 2536]

- หาวิธีในการตัดค่าที่สามารถจะเป็นไปได้ทั้งหมด
- เลือกแบบที่จำนวนค่าน้อยที่สุด
- เช่น เมื่อมีข้อความว่า "ไปห้ามเหสี" ก็จะตัดค่าได้ 2 เส้นทาง คือ
  - 1. ไป ห้าม เห สี      มีจำนวนค่าเท่ากับ 4
  - 2. ไป หาม เหสี      มีจำนวนค่าเท่ากับ 3
- เลือกเส้นทางที่ 2 เนื่องจากมีจำนวนค่าน้อยที่สุด



# ระยะแก้ไข (Edit Distance)

- ระยะแก้ไข ระหว่าง 2 สตริง คือจำนวนครั้งที่น้อยที่สุดที่ต้องการในการแก้ไข (เพิ่ม,ลบ,เปลี่ยน) สตริงหนึ่งให้เท่ากับอีกสตริงหนึ่ง
- วิทาลัย กับ วิทยาลัย ระยะแก้ไขเท่ากับ 1
- อาการ กับ อาหารรุ ระยะแก้ไขเท่ากับ 2

# การหาระยะแก้ไข

- Damerau-Levenshtein Dynamic Programming [Levenshtein, 1965]
  - Edit distance between string s and t
  - $\text{cost} = (s[i] == t[j])? 0 : 1$
  - $d[i, j] = \text{Min} (d[i-1, j]+1, d[i, j-1]+1, d[i-1, j-1]+\text{cost})$

	อ	า	ห	า	ร	
	0	1	2	3	4	5
อ	1	0	1	2	3	4
ห	2	1	1	1	2	3
า	3	2	1	2	1	2
ร	4	3	2	2	2	<b>1</b>

# การจับคู่แบบประมาณ (Approximate Matching)

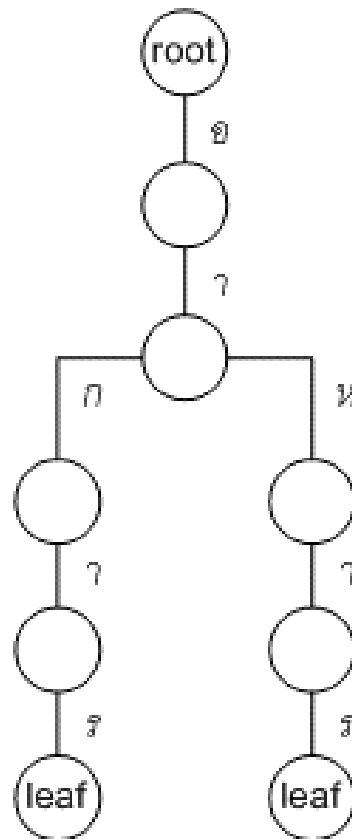
- การจับคู่ระหว่างคำซึ่งมีระยะแก้ไขไม่เกินค่าที่กำหนด
- ใช้ในการหาคำที่ต้องการจากเซตของคำที่เก็บไว้ (พจนานุกรม) โดยอนุญาตให้มีความผิดพลาดเกิดขึ้นได้ตามแต่จะกำหนด (Approximate Searching) เช่น
  - พจนานุกรมมีคำว่า “อาหาร” “อากร” “อาการ” “อาหฺร”
  - ค้นหาคำว่า “อาหฺร” โดยความผิดพลาดไม่เกิน 1 จะได้คำว่า “อาหฺร” “อาหาร” และ “อากร” เป็นคำตอบ

# ทรี (Trie)

- เป็นโครงสร้างข้อมูลแบบทรี (Tree) ซึ่งโหนดมีค่าว่าง (Empty)
- ข้อมูลจะอยู่บนเส้นทางการเดินระหว่างโหนดราก (Root Node) ไปยังโหนดใบ (Leaf Node)
- สามารถนำทรีมาใช้ในการเก็บคำในพจนานุกรมเพื่อค้นหาคำที่อยู่ในพจนานุกรม โดยแต่ละกิ่ง (Branch) ระหว่างโหนด จะมีอักขระกำกับอยู่

# ตัวอย่างทรีที่เก็บคำในพจนานุกรม

- เก็บคำว่า “อาการ” และ “อาหาร”



# ข้อดีของการใช้ทริยสำหรับเก็บคำในพจนานุกรม

- ใช้ส่วนหน้าของคำ (prefix) ร่วมกัน
- ประหยัดเนื้อที่ ง่ายต่อการเพิ่มคำ
- เวลาที่ใช้ไม่ขึ้นกับจำนวนคำในพจนานุกรม แต่ขึ้นกับความยาวของสตริงที่ต้องการค้นหา

# การค้นหาโดยการจับคู่แบบประมาณบนทรี

[H.Shang, T.H.Merrett, 1996]

- ค้นหาโดยเดินไปตามทรีแบบ Depth First Search, Recursive
- ขณะค้นหาเมื่ออยู่ที่ทรีย่อยใด ๆ พบความผิดพลาด หากพบ Ukkonen's Cutoff จะหยุดการค้นหาในทรีย่อยนั้น (ไม่จำเป็นต้องเดินให้ถึง leaf node) และย้อนกลับไปค้นหาในทรีย่อยถัดไป

# Ukkonen's Cutoff [Ukkonen, 1985]

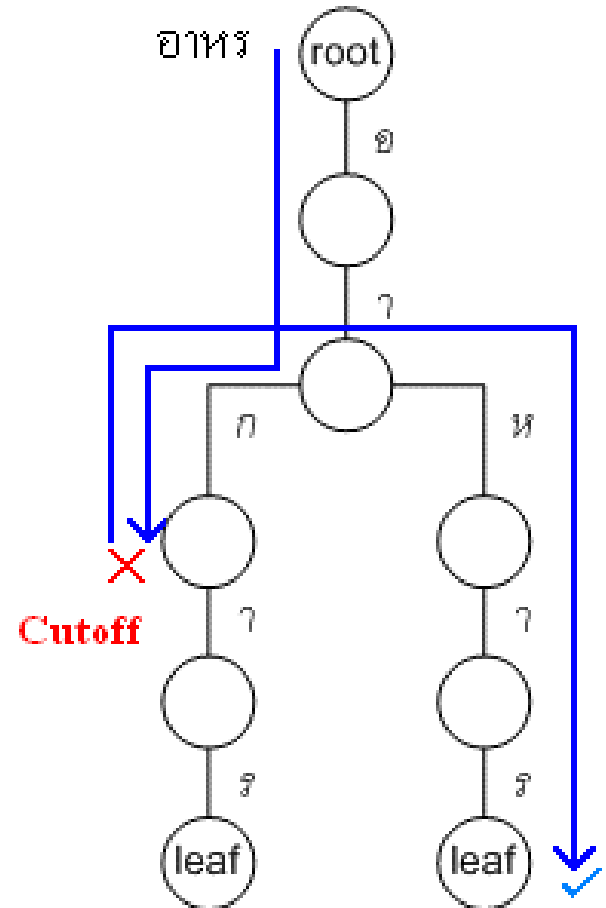
- ลดการคำนวณตาราง ในเซลล์ที่ไม่จำเป็นต้องคำนวณ เนื่องจากค่าที่นำมาคำนวณเกินความผิดพลาดที่ยอมรับได้แล้ว

	อ	๑	ก	๒	ร	
อ	0	1	2	3	4	5
๑	1	0	1	2		
ก	2	1	1	2		
๒	3	2	1	2		
ร	4	3	2	2		



# ตัวอย่างการค้นหาโดยการจับคู่แบบประมาณบนทรี

- ในพจนานุกรมมีคำว่า  
“อาการ” และ “อาหาร”
- คำที่นำค้นหาแบบประมาณ “อาหาร”



# แนวคิดการออกแบบอัลกอริธึมการตัดคำ

- จากวิธีการตัดคำแบบ Maximal Matching
- นอกจากพิจารณาจำนวนคำถูก ควรพิจารณาถึงโอกาสของคำที่จะเป็นคำผิด 1 อักขระด้วย
  - เดิม  $\# [\text{คำถูก} + \text{คำถูก} + \text{คำถูก} + \dots] > \# [\text{คำถูก}]$  ดังนั้น เลือกด้านขวา
  - ใหม่  $\# [\text{คำถูก} + \text{คำถูก} + \text{คำถูก} + \dots] > \# [\text{คำผิด 1 อักขระ}]$  ควรเลือกด้านขวา
- เช่น
  - $\# [\text{รัฐ} + \text{มน} + \text{ตรี}] > \# [\text{รัฐมนตรี}]$  ดังนั้นตัดคำเป็น รัฐมนตรี
  - $\# [\text{รัฐ} + \text{นม} + \text{ตรี}] > \# [\text{รัฐนมตรี}]$  ควรจะตัดเป็น รัฐนมตรี
  - $\# [\text{กระ} + \text{ทรง} + \text{วง}] > \# [\text{กระทรงวง}]$  ควรจะตัดเป็น กระทรงวง

## แนวคิดการออกแบบอัลกอริธึมการตัดคำ(2)

- การพิจารณาจากจำนวนคำเพียงอย่างเดียว อาจเกิดความผิดพลาดใน 2 กรณี
  1. คำถูก + คำถูก ไม่ควรตัดคำเป็น คำผิด
    - ว่า การ ไม่ควรตัดคำเป็น ว่าการ (วฏาการ)
    - จะ ต้อง ไม่ควรตัดคำเป็น จะต้อง (มะต้อง)
  2. คำถูก + คำถูก + คำถูก (ถูก 3) ไม่ควรตัดคำเป็น คำผิด + คำผิด (ผิด 2)
    - ทดลอง พบ ว่า -> ทดลองพ(ทดลอง) + บว่า(บ้า)

## แนวคิดการออกแบบอัลกอริธึมการตัดคำ (3)

- ดังนั้นควรกำหนดค่า cost ของคำถูกและคำผิด แล้วเลือกด้านที่ cost น้อยกว่า
- จากวิธีการเลือกตัดคำที่ผ่านมา ทำให้กำหนดความสัมพันธ์ของ cost ได้ดังนี้

กำหนดให้  $R = \text{cost of right word}$ ,  $W = \text{cost of misspelled word}$

- $\text{cost} [\text{คำถูก} + \text{คำถูก} + \dots] > \text{cost} [\text{คำถูก}] \quad (mR > nR; m > n) \text{ (Maximal Matching)}$
  - $\text{cost} [\text{คำถูก} + \text{คำถูก} + \text{คำถูก} + \dots] > \text{cost} [\text{คำผิด}] \quad (nR > W; n > 2) \Rightarrow 3R > W$
  - $\text{cost} [\text{คำถูก}] < \text{cost} [\text{คำผิด}]$
  - $\text{cost} [\text{คำถูก} + \text{คำถูก}] < \text{cost} [\text{คำผิด}] \quad (nR < W; n \leq 2) \Rightarrow 2R < W$
  - $\text{cost} [\text{คำถูก} + \text{คำถูก} + \text{คำถูก}] < \text{cost} [\text{คำผิด} + \text{คำผิด}] \quad (3R < 2W)$
- สรุปได้ว่า  $2R < W < 3R$
  - ต้องการทราบค่าของ  $W$  ที่เหมาะสม

# การทดสอบเพื่อหาค่าต้นทุนของคำผิดที่เหมาะสม

- Orchid Corpus ของ NECTEC
- สุ่มเลือกมา 1250 บรรทัด
- 250 บรรทัด สร้างความผิดพลาดบรรทัดละ 1 อักขระ ทั้ง 4 แบบ
- ถูกต้อง 1000 บรรทัด ผิดพลาด 1000 บรรทัด
- การวัดประสิทธิภาพ
  - $\text{Recall} = \frac{\text{จำนวนคำผิดที่พบ}}{\text{จำนวนคำผิดทั้งหมด}}$
  - $\text{Precision} = \frac{\text{จำนวนคำผิดที่พบและผิดจริง}}{\text{จำนวนคำผิดที่พบทั้งหมด}}$

## การทดสอบเพื่อหาค่าต้นทุนของคำผิดที่เหมาะสม(2)

ต้นทุนของคำผิด (W)	Recall	Precision
1.9R	$957 / 1000 = 95.7\%$	$957 / 1449 = 66.05\%$
2.0R	$839 / 1000 = 83.9\%$	$839 / 839 = 100\%$
2.5R	$839 / 1000 = 83.9\%$	$839 / 839 = 100\%$
2.9R	$839 / 1000 = 83.9\%$	$839 / 839 = 100\%$
3.0R	$806 / 1000 = 80.6\%$	$839 / 839 = 100\%$

# อัลกอริธึมการตัดคำสำหรับการตรวจสอบตัวสะกด

- ค้นหาคำที่เป็นไปได้ทั้งหมดในสตริง ดังนี้
  - คำที่พบในพจนานุกรม => คำถูก
  - คำที่ไม่พบในพจนานุกรม ผิดไม่เกิน 1 อักขระ ในพจนานุกรม => คำผิด
- เลือกเส้นทางตัดคำที่ต้นทุนน้อยที่สุด โดยพิจารณาดังนี้
  - คำถูก มีต้นทุนเท่ากับ 1
  - คำผิด มีต้นทุนเท่ากับ 2.5
  - $\text{ต้นทุน} = (\text{จำนวนคำถูก} * 1) + (\text{จำนวนคำผิด} * 2.5)$
  - หากต้นทุนเท่ากัน ใช้หลัก Longest Matching ในการเลือกเส้นทาง

# การทดสอบประสิทธิภาพ

	ขั้นตอนวิธีที่ นำเสนอ	ไมโครซอฟท์ เวิร์ด 2002	ปลาคาว ออฟฟิศ 2.0
จำนวนบรรทัดทั้งหมด	2000	2000	2000
จำนวนบรรทัดที่ถูก (AR)	1000	1000	1000
จำนวนบรรทัดที่ผิด (AW)	1000	1000	1000
คำที่ถูกระบุว่าผิด (IW)	839	757	783
คำผิดที่ถูกพบ (DW)	839	757	783
Recall (DW / AW)	<b>83.90%</b>	<b>75.70%</b>	<b>78.30%</b>
Precision (DW / IW)	100%	100%	100%
พบบนตำแหน่งที่ผิด (AP)	782	664	633
ความแม่นยำในการตรวจพบบนตำแหน่งที่ ผิด (AP / DW)	<b>93.21%</b>	<b>87.71%</b>	<b>80.84%</b>
แนะนำคำที่ถูกต้อง (AS)	713	320	258
ความสามารถในการแก้ไขเมื่อตรวจพบ (AS / DW)	<b>84.98%</b>	<b>42.27%</b>	<b>32.95%</b>



# วิเคราะห์ผลการทดสอบประสิทธิภาพ

- สามารถตรวจสอบคำผิดดีกว่าอีก 2 โปรแกรมเล็กน้อย
- สามารถในการแก้ไขเมื่อตรวจสอบพบได้ดีกว่ามาก
- เนื่องจากการตัดคำเมื่อข้อความมีความผิดพลาดทำได้ดีกว่าเมื่อเทียบกับอีก 2 โปรแกรม จึงทำให้สามารถระบุคำผิดได้ถูกต้อง และแนะนำคำที่ถูกต้อง

# ผลการทดสอบประสิทธิภาพ แบ่งตามประเภทของความผิดพลาด

ประเภทความ ผิดพลาด	ทั้งหมด	ตรวจสอบ พบ	ความสามารถในการ ตรวจพบ
พิมพ์เกิน	250	221	88.4%
พิมพ์ผิด	250	207	82.8%
พิมพ์ตก	250	177	70.8%
พิมพ์สลับ	250	234	93.6%
รวม	1000	839	83.9%

## วิเคราะห์ผลการทดสอบประสิทธิภาพ (2)

- พิมพ์ตกจะสามารถตรวจสอบพบได้น้อยที่สุด เนื่องจากคำในพจนานุกรมส่วนใหญ่จะเป็นคำเล็ก การพิมพ์ตกทำให้อักขระลดน้อยลง จึงมีโอกาที่จะเป็นคำถูกแม้ว่าจะพิมพ์ตก จึงตรวจสอบพบได้ยากกว่าแบบอื่น ๆ
- พิมพ์สลับสามารถตรวจสอบพบได้มากที่สุด เพราะอาจมองเป็นการพิมพ์ผิดถึง 2 ที่ จึงมีโอกาสน้อยที่จะยังคงเป็นคำถูกอยู่หลังจากการสลับที่ จึงตรวจสอบพบได้ง่ายกว่าแบบอื่น ๆ

# การทดสอบความเร็วของขั้นตอนวิธี

รหัสกรณีทดสอบ	กรณีทดสอบ	จำนวนอักขระ	เวลาที่ใช้ (วินาที)
TDS6-001	การประชุมทางวิชาการ โครงการวิจัยและพัฒนา ปีงบประมาณศูนย์เทคโนโลยีอิเล็กทรอนิกส์และ คอมพิวเตอร์แห่งชาติศูนย์เทคโนโลยีอิเล็กทรอนิกส์และ คอมพิวเตอร์แห่งชาติ	148	8.922
TDS6-002	รัฐมนตรีว่าการกระทรวงวิทยาศาสตร์ประเทศไทยได้มีการ ปรับเปลี่ยน โครงสร้างในการพัฒนาเศรษฐกิจของประเทศ จากประเทศเกษตรกรรมไปสู่ความเป็นประเทศอุตสาหกรรม	141	8.625
TDS6-003	มนุษย์ส่วนใหญ่ในสังคมจะต้องประสบกับปัญหาสังคมไม่ กระทบกระเทือนไม่กระทบกระเทือนเพื่อให้ปัญหานั้น บรรเทาเบาบางลงการจรรยาบรรณจิตขัดความจริงไม่ได้เป็นหลัก ประกันว่าจะแก้ปัญหาได้	164	9.890
TDS6-004	ความล้มเหลวของกฎเกณฑ์ที่สังคมกำหนดแต่เป็นเพียง สภาพการณ์ที่ไม่สอดคล้องกับคุณค่าที่กลุ่มยึดถือแต่เป็นเพียง สภาพการณ์ที่ไม่สอดคล้องกับคุณค่าที่กลุ่มยึดถือใน ระบบประมวลกฎหมาย	164	9.312

# ข้อจำกัด

- ไม่สามารถตรวจสอบคำผิดที่สามารถแยกย่อยเป็นคำถูก 2 คำ
  - สอดคล้อง สะกดผิดเป็น สอดคล้อง จะตัดคำได้ สอดคล้อง  
(ที่ถูกต้องควรตัดเป็น สอดคล้อง)
- การสลับอักขระระหว่างคำ หรือระหว่างประโยค ซึ่งทำให้ส่วนที่เหลือของคำใดคำหนึ่งถูก
  - สูงมาก สะกดผิดเป็น สูงมาก จะตัดคำได้ สูงมาก  
(ที่ถูกต้องควรตัดเป็น สูงมาก)  
(“สูง” เป็นคำถูก ทำให้ “มาก” ถูกมองเป็นการพิมพ์เกิน)

## ข้อจำกัด(2)

- คำที่มี 2 อักขระ หากพิมพ์ตกไป 1 อักขระ อักขระที่เหลือจะถูกมองเป็นการพิมพ์เกินของคำข้างเคียง
  - ต่อไป สะกดผิดเป็น ต่อป จะตัดคำได้ ต่อป  
(ที่ถูกตัดควรตัดเป็น ต่อ|ป)
- กรณีที่ต้นทุนเท่ากัน แล้วใช้วิธี Longest Matching ในการเลือกตัดคำ อาจ ได้ผลลัพธ์ที่ไม่ถูกต้อง
  - วิธีแสดง สะกดผิดเป็น วิธีสแดง จะตัดคำได้ วิธีส|แดง  
(ที่ถูกตัดควรตัดเป็น วิธี|สแดง)

# ข้อเสนอแนะ

- ขั้นตอนวิธีจับคู่ประมาณสำหรับการค้นหาคำผิดในสตริงหนึ่ง ๆ ควรได้รับการปรับปรุงให้เร็วขึ้น
- การตรวจสอบตัวสะกดในกรณีที่เป็นคำถูกอยู่ติดกัน 2 คำ ซึ่งเป็นข้อจำกัดของอัลกอริทึมนี้ สามารถใช้ข้อมูลทางสถิติ หรือทางไวยากรณ์ บ่งบอกคำผิดได้

# จบการนำเสนอ