

บุญวัฒน์ ธาดาภักย์ 2554: การระบุเว็บโฮสต์ไทยโดยใช้เครื่องจักรเรียนรู้ ปริญญา
วิศวกรรมศาสตรมหาบัณฑิต (วิศวกรรมคอมพิวเตอร์) สาขาวิศวกรรมคอมพิวเตอร์
ภาควิชาวิศวกรรมคอมพิวเตอร์ อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก: ผู้ช่วยศาสตราจารย์
อานนท์ รุ่งสว่าง, Ph.D. 79 หน้า

เว็บเพจบนอินเทอร์เน็ตถือว่าเป็นแหล่งความรู้ขนาดใหญ่ และมีการเพิ่มจำนวนขึ้น
อย่างต่อเนื่อง จนทำให้ไม่สามารถเก็บรวบรวมเว็บเพจทั้งหมดได้ภายใต้ทรัพยากรระบบที่มีอยู่
อย่างจำกัด เช่น ระยะเวลา พื้นที่เก็บข้อมูล หรือแบนด์วิดท์ เป็นต้น ด้วยเหตุนี้การใช้เว็บคราเวลอร์
แบบเจาะจงภาษาจึงเป็นทางเลือกหนึ่งที่จะมาช่วยให้สามารถเลือกเก็บเฉพาะเว็บเพจที่เขียนด้วย
ภาษาตามที่ต้องการได้ ตัวอย่างการนำไปใช้งานเช่น ผู้ให้บริการระบบสืบค้นในแต่ละประเทศมักจะ
แสดงเว็บเพจผลลัพธ์ตามภาษาประจำชาติของผู้ใช้ เพื่อให้ผู้ใช้เกิดความพึงพอใจต่อระบบ หรือ
การสร้างคลังเว็บเพจประจำชาติ เพื่อเก็บรักษาเว็บเพจที่มีข้อมูลสำคัญด้านสังคม วัฒนธรรม
วิถีการดำเนินชีวิตให้กับคนรุ่นต่อไป

วิทยานิพนธ์ฉบับนี้นำเสนอระบบต้นแบบเว็บคราเวลอร์เจาะจงเว็บเพจภาษาไทย
แบบแยกเก็บตามไซต์ ซึ่งประกอบด้วยส่วนของการทำนายหรือระบุภาษาของ โฮสต์ และส่วนของ
เว็บคราเวลอร์แบบแยกเก็บตามไซต์ โดยที่ส่วนของการทำนายภาษาใช้เทคนิคทางเครื่องจักรเรียนรู้
แบบหลายตัว ซึ่งเครื่องจักรเรียนรู้แต่ละตัวจะเรียนรู้จากคุณลักษณะทั้งตัวอย่างของเว็บโฮสต์
ที่เกี่ยวข้อง และเว็บโฮสต์ที่ไม่เกี่ยวข้อง ด้วยเทคนิคทางเครื่องจักรเรียนรู้ที่แตกต่างกัน และ
ส่วนทำนายภาษายังมีหน้าที่ช่วยระบุภาษาของเว็บโฮสต์เป้าหมายพร้อมทั้งให้คะแนนความมั่นใจ
อีกด้วย แต่สำหรับเว็บคราเวลอร์แบบแยกเก็บตามไซต์มีหน้าที่เก็บเว็บเพจภายใต้เว็บไซต์ที่กำหนด
เรียงลำดับตามคะแนนความมั่นใจที่ได้จากส่วนการทำนายภาษา จากผลการทดลองแสดงให้เห็นว่า
ประสิทธิภาพการเก็บเว็บเพจจากอินเทอร์เน็ตของวิธีการที่นำเสนอดีกว่าวิธีการของเว็บคราเวลอร์
เจาะจงภาษาในระดับเว็บเพจที่ผ่านมา ซึ่งวัดผลจากอัตราการเก็บเกี่ยวเว็บเพจภาษาไทย

ลายมือชื่อนิติสด

ลายมือชื่ออาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก