

Project Proposal

ชื่อโครงการ: การเปรียบเทียบPageIndex และ PageIndex ร่วมกับ Light Knowledge Graph สำหรับระบบ
查詢บทอทตอบคำถมภมายอาญไทย

1. ปัญหาและความสำคัญ (Problem Statement)

ระบบ Retrieval-Augmented Generation (RAG) เป็นเทคนิคสำคัญในการสร้าง AI ที่สามารถตอบ
คำถมจากเอกสาร โดยระบบ RAG แบบดั้งเดิมใช้ Vector Embeddings ในการค้นหาข้อความที่มีความคล้ายคลึง
ทางความหมาย อย่างไรก็ตาม วิธีนี้มีข้อจำกัดสำคัญสำหรับเอกสารภมาย ได้แก่ (1) ข้อความที่มีความคล้ายคลึง
กันแต่ไม่ได้สอดคล้องกันทางความหมาย เช่น คำถมเรื่อง "โทษลักษทรพย" อาจจึงข้อความเรื่อง "โทษชิงทรพย" มา
แทน (2) การแบ่งเนื้อหาเอกสารออกเป็นส่วนที่เท่ากันทำลายโครงสร้างลำดับขั้นของภมาย และ (3) การค้นหา
แบบ Vector ไม่สามารถให้เหตุผลหลายขั้นตอนสำหรับคำถมภมายที่ซับซ้อนได้

PageIndex จึงเป็นแนวทาง RAG แบบที่ไม่ต้องใช้ Vector ที่เสนอขึ้นใหม่เพื่อแก้ปัญหาดังกล่าว โดยการ
สร้างดัชนีแบบต้นไม้จากเอกสารและใช้ LLM เข้ามาช่วยในการให้เหตุผลเพื่อค้นหา อย่างไรก็ตาม ในปัจจุบันยัง
ไม่มีการศึกษาว่า PageIndex จะทำงานได้ดีเพียงใดกับเอกสารภมายภาษาไทย และยังไม่มีงานวิจัยที่ศึกษา
การนำ Knowledge Graph มาช่วยเพิ่มความแม่นยำของระบบ

2. วัตถุประสงค์ (Objectives)

2.1 พัฒนาระบบ PageIndex สำหรับประมวลภมายอาญไทย พร้อมรองรับภาษาไทย

2.2 สร้าง Light Knowledge Graph ที่ประกอบด้วย entities ทางภมาย เช่น มาตรา, ความผิด, โทษ
และความสัมพันธ์

2.3 เปรียบเทียบประสิทธิภาพระหว่าง PageIndex อย่างเดียว กับ PageIndex ร่วมกับ Light
Knowledge Graph ในงาน查詢บทอทตอบคำถมภมายไทย

3. ชุดข้อมูล (Dataset)

3.1 คลังข้อมูลภมายไทยประเกพราชบัญญัติ ซึ่งรวบรวมมาจากเว็บไซต์ของสำนักงานคณะกรรมการ
การกฎหมาย (<https://www.krisdika.go.th/>) และที่มา <https://github.com/PyThaiNLP/thai-law?tab=readme-ov-file>

3.2 ชุดข้อมูลทดสอบ Q&A ที่สร้างขึ้น พร้อมคำตอบมาตรฐานและการอ้างอิงมาตรา

4. เทคนิค Machine Learning ที่ใช้ (Proposed ML Techniques)

| องค์ประกอบ | เทคนิค | วัตถุประสงค์ |
|---------------------|-------------------------------|------------------------------------------------|
| Document Indexing | PageIndex Tree Construction | สร้างโครงสร้างต้นไม้ตามลำดับชั้นของ |
| Text | LLM-based Summarization | สร้างสรุปสำหรับแต่ละ node ในต้นไม้ |
| Entity Extraction | Named Entity Recognition ผ่าน | ดึง entities ทางกฎหมาย (มาตรา, ความผิด, |
| Relation Extraction | LLM-based Relation Extraction | ระบุความสัมพันธ์ระหว่าง entities |
| Query | Intent Classification | จำแนกประเภทคำถามเพื่อเลือกกลยุทธ์การ |
| Retrieval | Tree Search + KG Lookup | การค้นหาแบบสมมผาน tree navigation กับ KG index |
| Answer Generation | RAG with LLM | สร้างคำตอบพร้อมอ้างอิงมาตรากฎหมาย |
| Evaluation | LLM-as-Judge, ROUGE, | ประเมินคุณภาพและความถูกต้องของคำตอบ |

5. ความท้าทายที่คาดว่าจะพบ (Expected Challenges)

- 1) การประมวลผลภาษาไทยและศัพท์กฎหมายเฉพาะทาง
- 2) การแยกวิเคราะห์โครงสร้างมาตรากฎหมาย
- 3) ความแม่นยำในการดึง Entity จาก LLM
- 4) ข้อจำกัดของชุดข้อมูลทดสอบขนาดเล็ก
- 5) การกำหนด Ground Truth สำหรับคำตามกฎหมายที่อาจมีคำตอบถูกต้องได้หลายแบบ