

RETRIVAL AUGMENTED GENERATION

András Kornai

2024 October 16

WHAT WE PLAN ON DOING

- Build a testbed
- Build a full system
- Do something for Hungarian
- Understand the theory better

READINGS

- **Today:** Church, Sun Yue, Vickers, Saba, Chandrasekar 2024: Emerging trends: a gentle introduction to RAG
- Still looking for takers: Zhao and Rios 2024, https://huggingface.co/learn/cookbook/en/rag_evaluation

SEPARATE TASKS: 1. INFRASTRUCTURE

- Download Hungarian WP **DONE** (Martin Juhász)
- Parse out tables and infoboxes **unassigned**
- Create a big database with lots of tables **unassigned**
- Create NL descriptions of tables 'Ez a táblázat Magyarország miniszterelnökeit listázza, dátumokkal' **unassigned**
- Verify DB purity (make sure stuff is not in other WPs) **unassigned**
- Create join chains "X was minister of Y government, Z is son of X" ... **unassigned**

SEPARATE TASKS: 2. LLM

- Find a range of pre-existing LLMs
- Compare them on standard (Hungarian) tasks
- Maybe train adaptation layer
- Maybe add chain-of-thought prompting
- Test how these work

SEPARATE TASKS: 3. RETRIEVAL

- Find appropriate docs outside WP
- Segment WP
- Construct retrieval system based on keywords – Lucene
- Construct retrieval system based on triple storage – if Istvan Soltesz finishes his port of OpenIE for Hungarian
- Construct retrieval system based on vectors – Milvus
- Maybe combine these three, evaluate on its own merits

BY OCT 30: BUILD AND TEST A BASELINE SYSTEM

- Trello has the individual responsibilities
- We will have a zoom only meeting on the 23rd at 4PM (usual zoom link)
- We will use HuBERT for generating the vectors for the snippets
- Watch slack, we may have to move off of trello (monetization)