

Example Slide (~5min) - due 2/6

- Presenter: 한상은
- Topic: Grounded QA
- Data and Baseline:
[https://github.com/Snowflake-Labs/Arctic_Agentic_RAG? fsi=urgDB5P9](https://github.com/Snowflake-Labs/Arctic_Agentic_RAG?_fsi=urgDB5P9)
- Novel Hypothesis: This model may not perform well when generator is weaker. **We will propose** new model consistently working well for weak/strong generator (you can learn from negative results too— just **sharing in academic form** is what we aim for)
- Output: github repo, ppt

Final Presentation (~7min) - due 2/10

- Example Final Slide
 - https://drive.google.com/file/d/1f_RCptUpQqVwhb8elfkz7g2P6KhB_Gqr/view?usp=sharing
 - You can make it simpler
- Motivation, Method, Experiment

Github Repo

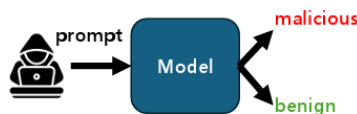
- Student project from last semester:
 - <https://github.com/Superfish83/PromptKiller/>
- Provide access to your dataset/model, explain how to run the code

Proposal of Idea & Contribution

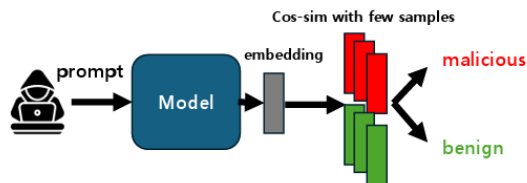
Idea: Few-shot learning will enable detectors to filter out novel malicious prompts, of which we have very few samples.

-> **Siamese fine-tuning with prompt data**

Existing classification approaches



My model: Few-shot (Siamese) approach



Github Repo

- Tutorial for creating/updating github repo
 - <https://docs.github.com/en/repositories/creating-and-managing-repositories/quickstart-for-repositories>
- Walk-through demo
 - <https://www.youtube.com/watch?v=-9rTsx6ilvo>
- Evaluation (30)
 - Creativity (독창성): 새로운 가설 및 베이스라인 개선방법을 제시하였는가?
 - Relevance (업무 연관성): NLP 기술을 업무에 접목시켰는가? (e.g. 업무관련 benchmark 이용)
 - Depth (깊이):
 - 제시하는 모델과 방법론을 왜 선택했는지 기술적으로 설명할 수 있는가?
 - 다루는 문제에 대응되는 적절한 방법(e.g. Data scarcity 문제를 해결하기 위한 Data augmentation 기법)이 사용되었는가?

Update

- Topic: Synthetic, hypothetical dataset is OK (llm으로 생성 가능)
- Baseline: Vanilla LLM prompting
- Method: Customized agentic workflow (실습에서 배운 LangChain을 이용)
- 금요일 발표: topic, dataset, agentic workflow 도입 계획 (~5min)
 - Not graded
- 화요일 발표: experimental result / analysis 포함 (~7min)
 - Graded