# Research Progress on Data Intelligence in AI Graduate Program (AI 대학원 데이터 지능 연구 성과)

**Jongwuk Lee (이종욱)**
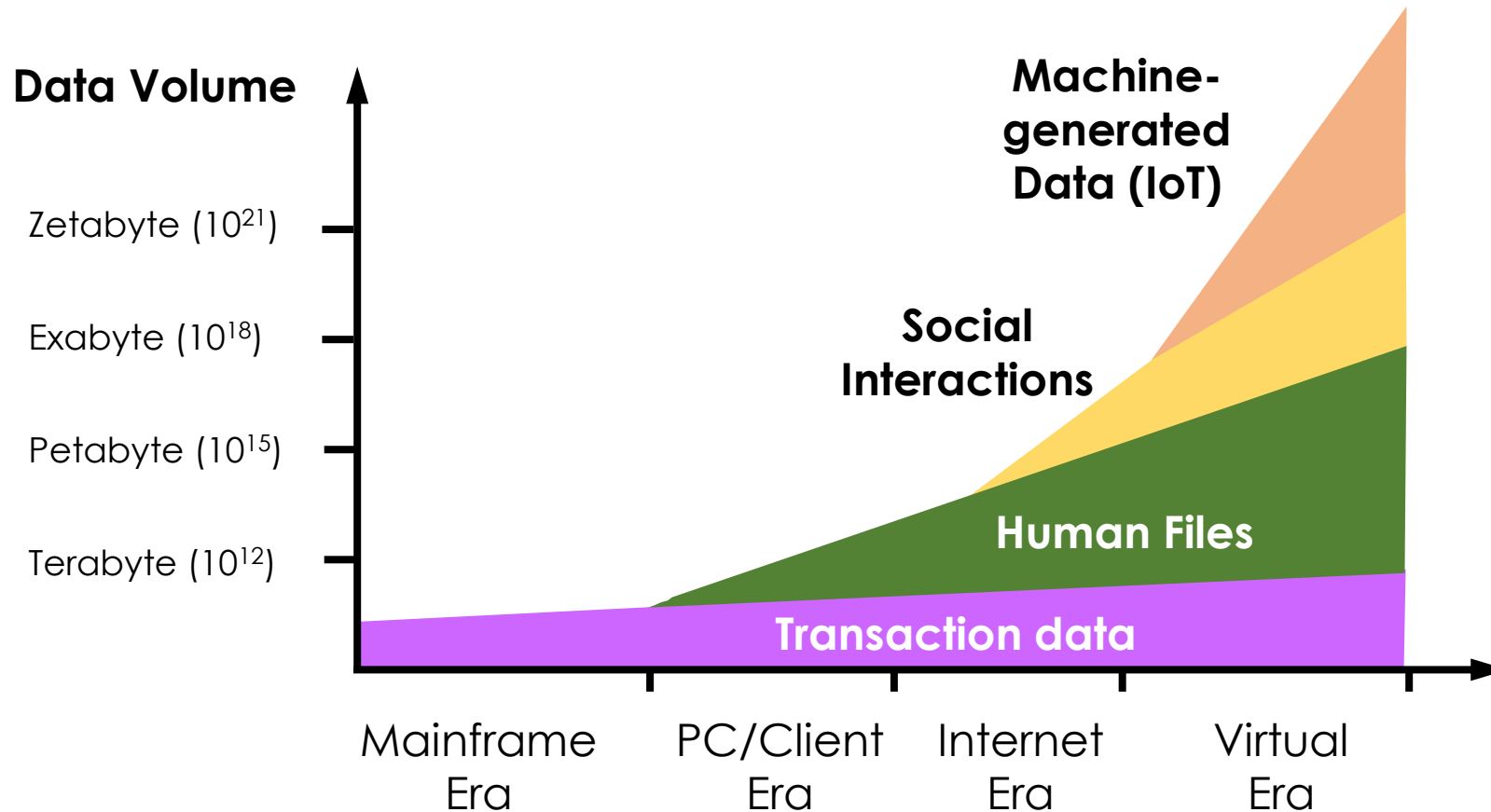
**Sungkyunkwan University (성균관대)**

# What is Data Intelligence?

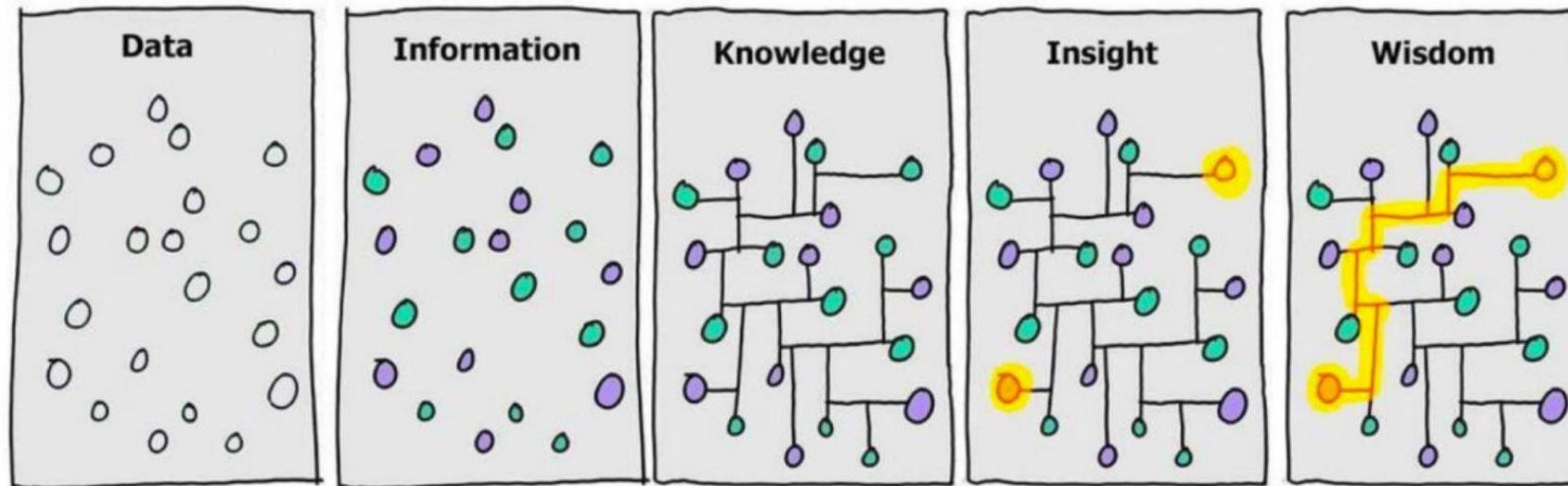# Explosive Growth of Data

- **Data exist everywhere!**

# Data Flooding and Overload

- **Drowning in data** but **starving for knowledge**!

# What is Data Mining (DM)?

- **Data mining (knowledge discovery from data)**
  - Extraction of **interesting** (**non-trivial**, **implicit**, **unknown** and **potentially useful**) **patterns** or **knowledge** from data

- **Refers to data intelligence for business/real-world objective.**

# Example: Word Association Map

- **Finding relevant words for "코로나19" from a news corpus**



2020년 11월 4주차
(2020.11.23 ~ 2020.11.25)

"코로나19" 감성 연관어 TOP 10

힘들다 · 바라다 · 성공하다 · 우려 · 확산 · 안전 · 크다 · 적극적 · 위기 · 저렴하다

긍정 45%  부정 23%  중립 32%

감성어 랭킹

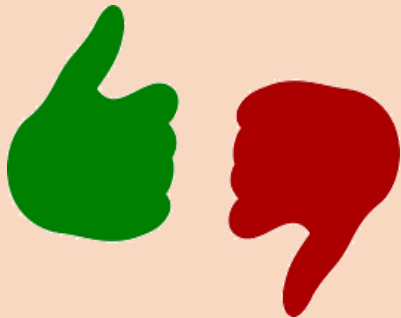| 순위 | 분류 | 키워드 | 건수 |
|---|---|---|---|
| 1 | 중립 | 확산 | 6,965 |
| 2 | 긍정 | 안전 | 5,059 |
| 3 | 부정 | 위기 | 2,627 |
| 4 | 긍정 | 성공하다 | 1,957 |
| 5 | 부정 | 힘들다 | 1,801 |
| 6 | 긍정 | 적극적 | 1,663 |
| 7 | 부정 | 우려 | 1,628 |
| 8 | 긍정 | 저렴하다 | 1,491 |
| 9 | 긍정 | 바라다 | 1,423 |
| 10 | 중립 | 크다 | 1,368 |

# Our Achievements in AI Graduate Program

- **20+ top-tier papers** were published in WWW, KDD, SIGIR, ICDM, CIKM, ACL, EMNLP, COLING, and VLDB.

# Our Achievements in AI Graduate Program

- **Three main research topics are**

**Rec systems**

**Graph mining**

**Text mining and Understanding**

# Recommender Systems

# Case 1: Recommender Systems

- **How to provide relevant items to users?**

Systems take initiative. (**push mode**)

**Recommendations**

**Items**

**products, movies, music, news, …**

# Our Achievements in RecSys

- **6 papers in WWW, SIGIR, ICDM, and CIKM**
  - Dual Neural Personalized Ranking, **WWW** 2019 (**SKKU**)
  - Collaborative Distillation for Top-N Recommendation, **ICDM** 2019 (**SKKU**)
  - AR-CF: Augmenting Virtual Users and Items in Collaborative Filtering for Addressing Cold-Start Problem, **SIGIR** 2020 (**Hanyang Univ.**)
  - Interest Sustainability-Aware Recommender System, **ICDM** 2020 (**POSTECH**)
  - DE-RRD: A Knowledge Distillation Framework for Recommender System, **CIKM** 2020 (**POSTECH**)
  - News Recommendation with Topic-Enriched Knowledge Graphs, **CIKM** 2020 (**Yonsei Univ.**)

# Dual Pairwise Ranking

- **Utilizing both <span style="color:red">user- and item-side pairwise rankings</span> over a neural architecture**

Utilizing **dual pairwise rankings**

Pos. Items **User** Neg. Items

Pos. Users **Item** Neg. Users

Incorporating one **generalized DNNs**

- **Adopting GANs to address the cold-start problem**



**Training CGANs to generate virtual users/items**

**Augmenting the original user-item matrix**

# Knowledge Distillation (KD)

- A small **student model** is trained to mimic a pre-trained and large **teacher model**.

- **Applying KD for recommender models**

*Collaborative Distillation for Top-N Recommendation, ICDM 2019 (SKKU)*

# Distillation Experts for Compression

- **Propose distillation experts (DE) to transfer latent knowledge from the teacher model.**

# Graph Mining

# Case 2: Graph Analysis and Learning

- **Finding frequent subgraphs and substructures**
- **Learning graph embedding**
- **Predicting/classifying node and edges**



**5B+ Web pages**

**6k+ proteins**

**Knowledge graph**

# Case 2: Graph Analysis and Learning

- **Nodes as people and edges as friendship**



**2.5B+ users in Facebook**

# Our Achievements in Graph Mining

- **8 papers in WWW, KDD, SIGIR, ICDM, and VLDB**
  - How Much and When Do We Need Higher-order Information in Hypergraphs? A Case Study on Hyperedge Prediction (**WWW** 2020, **KAIST**)
  - Structural Patterns and Generative Models of Real-world Hypergraphs (**KDD** 2020, **KAIST**)
  - SSumM: Sparse Summarization of Massive Graphs (**KDD** 2020, **KAIST**)
  - Incremental Lossless Graph Summarization (**KDD** 2020, **KAIST**)
  - Hypergraph Motifs: Concepts, Algorithms, and Discoveries (**VLDB** 2020, **KAIST**)
  - Evolution of Real-world Hypergraphs: Patterns and Models without Oracles (**ICDM** 2020, **KAIST**)
  - Unsupervised Differentiable Multi-aspect Network Embedding (**KDD** 2020, **POSTECH**)
  - ASiNE: Adversarial Signed Network Embedding (**SIGIR** 2020, **Hanyang Univ.**)

- **Generating a lossless graph summarization**

*The initial input has 9 nodes with 10 edges.*



Delete {f, i}

Add {f, i}

Delete {f, i}

Add {f, i} & Delete {f, i}

*The final output has 3 nodes with 4 edges.*

# Lossless Graph Summarization

- **How to effectively generate a lossless graph summarization for a dynamic environment?**
  - When the graphs are dynamically changed, minimizes the updates without reconstruction.

- **How to efficiently summarize the graph with a given bit constraint?**
  - **Given**: a graph $G$
  - **Find**: a summary graph $\overline{G}$
  - **Objective**: Minimize the difference between $G$ and $\overline{G}$.
  - **Subject to**: the size of $\overline{G}$ in bits $\leq k$.

*Incremental Lossless Graph Summarization, KDD 2020 (KAIST)*

*SSumM: Sparse Summarization of Massive Graphs, KDD 2020 (KAIST)*

# Graph Embedding

- **The similar nodes in a network have similar vector representation.**
  - The similarity in the embedding space **approximates** the similarity in the original network.

**Parameters initialized randomly**

**2-dim output per node**

[Karate club network]

*Source: https://iamsiva11.github.io/graph-embeddings-2017-part1/*

# Multi-aspect Graph Embedding



Context: "**schoolmate**"

Assign "**schoolmate**" aspect

Effective for **capturing multi-aspect user behavior**

Considering **multi-hop neighbors**

*Unsupervised Differentiable Multi-aspect Network Embedding, KDD 2020 (POSTECH)*

# Signed Graph Embedding

- **How to represent nodes in the embedding space?**
  - Nodes with the positive edges to be close
  - Nodes with the negative edges to be distant



- **A theory for signed graphs**

> "A friend (+) of my friend (+) is my friend (+)"
> "A friend (+) of my enemy (-) is my enemy (-)"
> "An enemy (-) of my friend (+) is my enemy (-)"
> "An enemy (-) of my enemy (-) is my friend (+)"

# Text Mining and Understanding

# Case 3: Text Mining and Understanding

- **Various NLP applications by understanding text**
  - Lage-scale text classification
  - Reading comprehension for question answering
  - Translating human language to machine language



**Donald Trump criticizes Dodgers manager for bullpen moves.**

# Case 3: Text Mining and Understanding

- **Natural language processing (NLP) lets computers to process and analyze large accounts of natural language data.**
  - Human-computer interaction
  - Includes the automation of linguistic forms, activities, and methods of communication.

**AI** + **Linguistics** = **NLP**

# Our Achievements in Text Mining

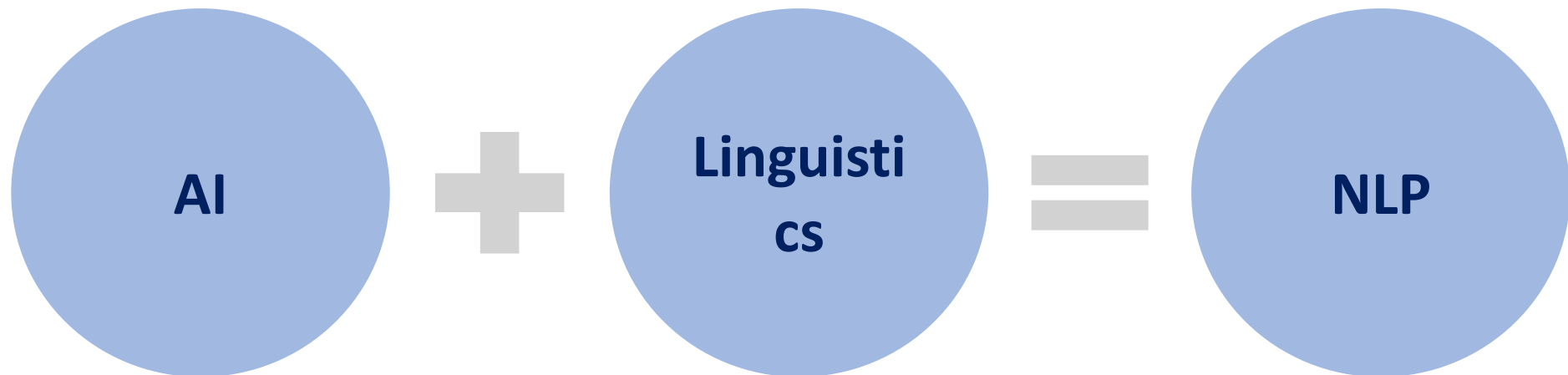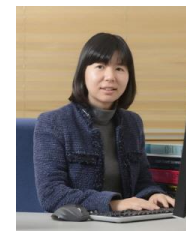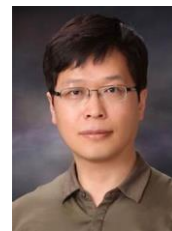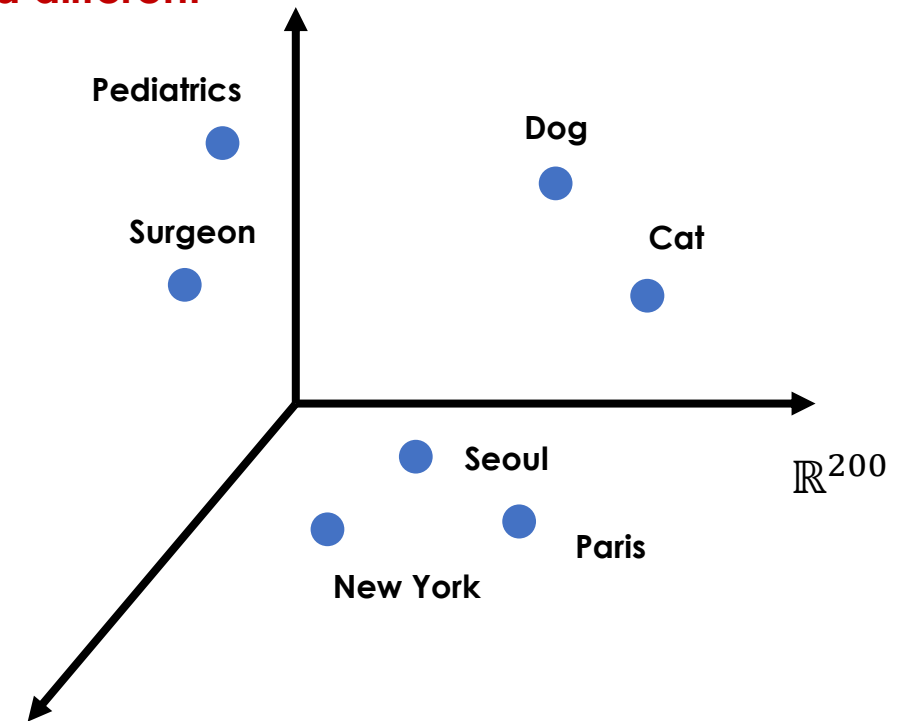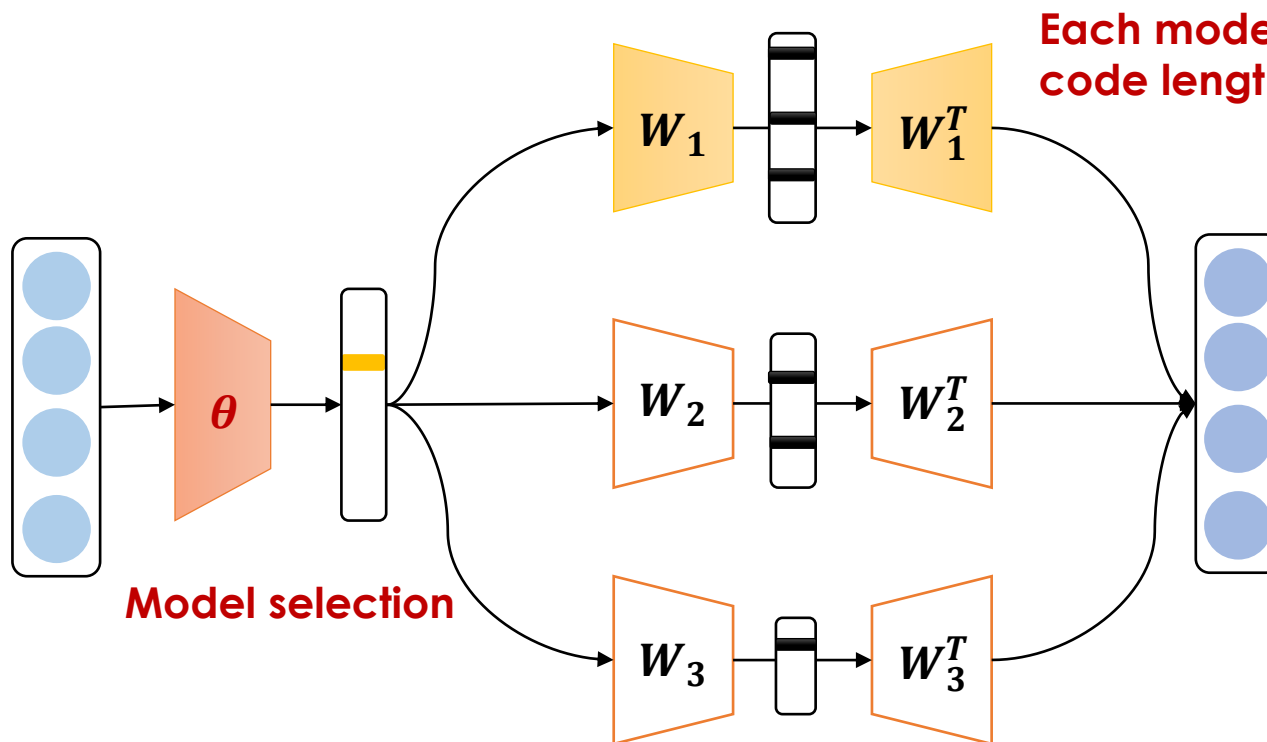- ## 9 papers in ACL, EMNLP, COLING, CIKM, and VLDB
  - Adaptive Compression of Word Embeddings, **ACL** 2020 (**Korea Univ.**)
  - Multi-pretraining for Large-scale Text Classification, **EMNLP** 2020 (**Korea Univ.**)
  - ST-GRAT: A Novel Spatio-temporal Graph Attention Network for Accurately Forecasting Dynamically Changing Road Speed, CIKM 2020 (**KAIST**)
  - Multi-Task Learning for Knowledge Graph Completion with Pre-trained Language Models, **COLING** 2020 (**SKKU**)
  - Natural language to SQL: Where are we today? **VLDB** 2020, (**POSTECH**)
  - Extracting Chemical–Protein Interactions via Calibrated Deep Neural Network and Self-training, **EMNLP** 2020 (**GIST**)
  - Learning with Limited data for Multilingual Reading Comprehension, **EMNLP** 2019 (**Yonsei Univ.**)
  - Less is More: Attention Supervision with Counterfactuals for Text Classification, **EMNLP** 2020 (**Yonsei Univ.**)
  - Retriever-Augmented and Controllable Review Generation, **COLING** 2020 (**Yonsei Univ.**)

# Word Embeddings Compression

- **Word embedding compression by adaptively assigning different lengths of discrete codes**
  - Using codes with a longer length for task-sensitive words



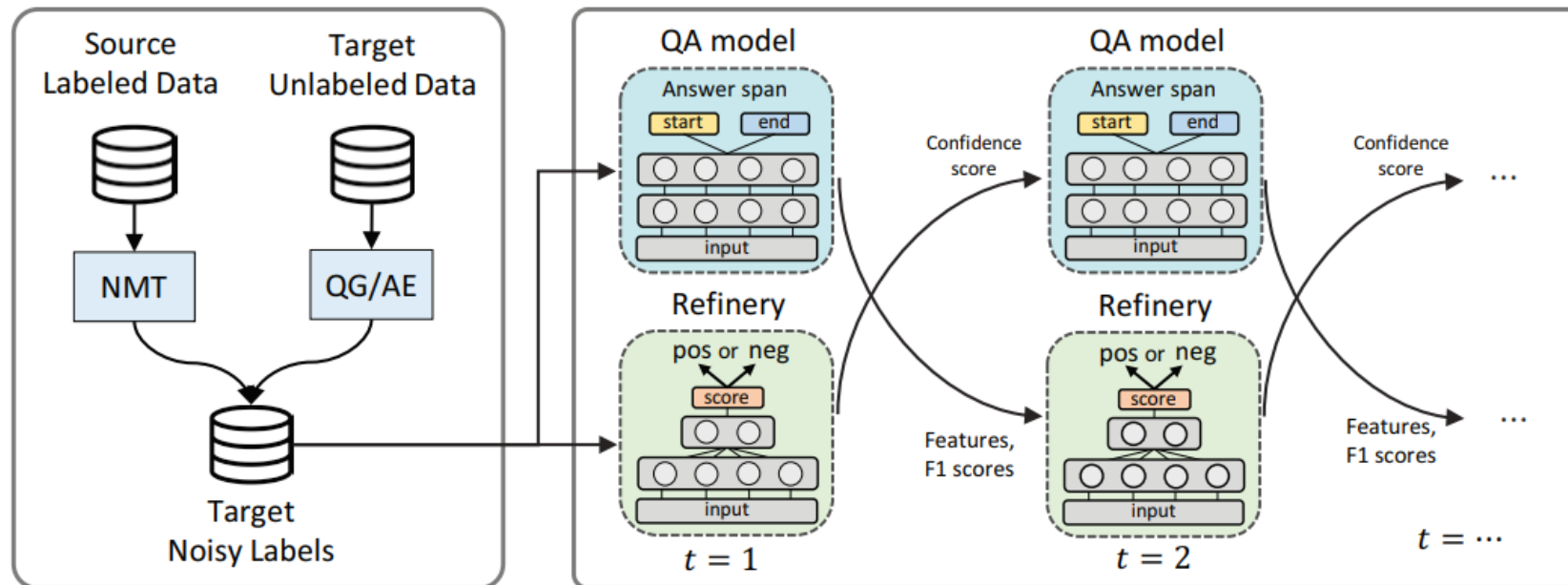Each model has a different code length.

Model selection
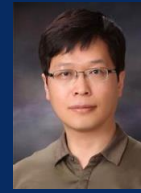
Example of word embeddings
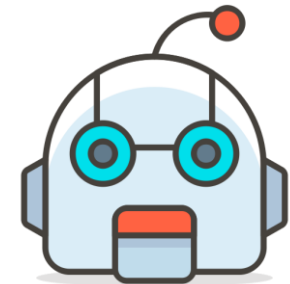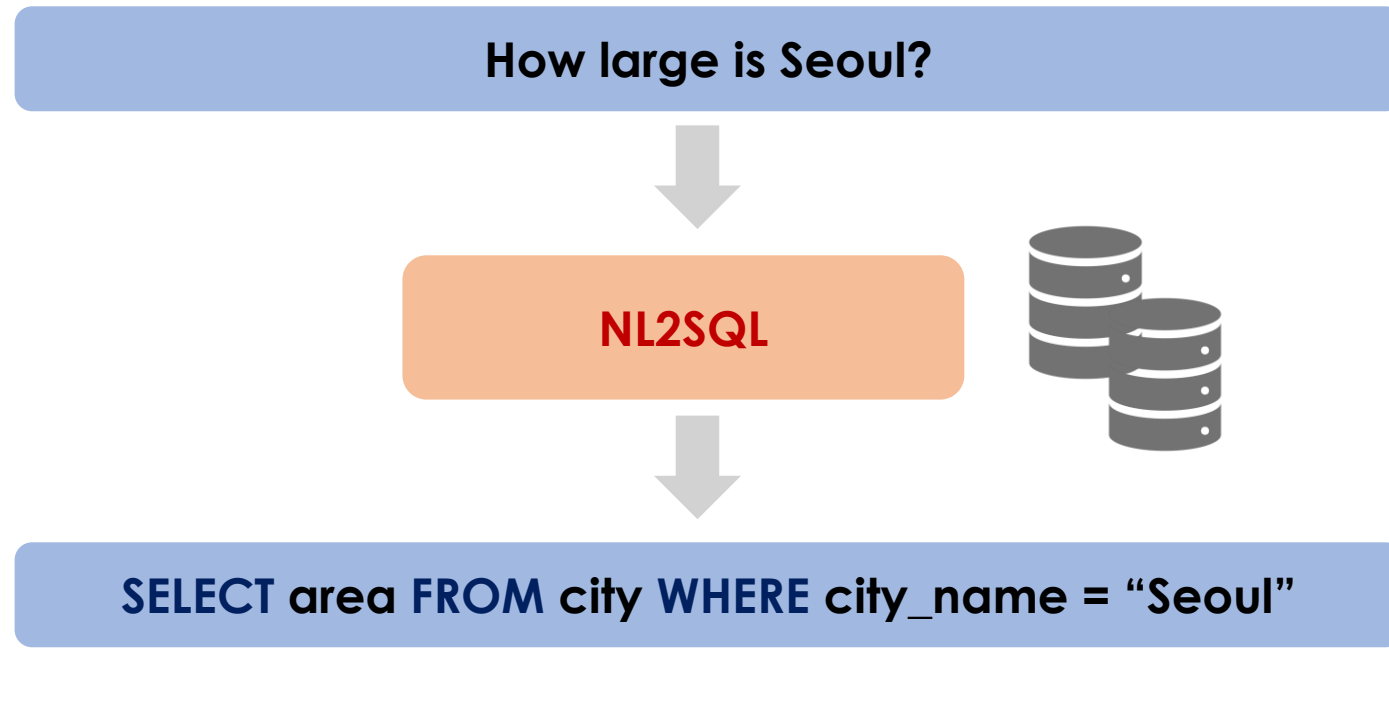
# Multilingual Reading Comprehension

- **Supporting question answering in a new language with limited training resources**
  - **Transfer labels** from another language.
  - Generate labels from unlabeled data using a **translator** and an **automatic labeling function**.

# NLP-to-SQL Translation

- **Generating an SQL statement to answer a natural language question on a relational database**

How large is Seoul?

NL2SQL

SELECT area FROM city WHERE city_name = "Seoul"

- **Provide a comprehensive survey for NL2SQL.**

*Natural language to SQL: Where are we today? VLDB 2020 (POSTECH)*

# Road Traffic Speed Prediction

- **Self-attentional model for forecasting spatio-temporal data**

*ST-GRAT: A Novel Spatio-temporal Graph Attention Network for Accurately Forecasting Dynamically Changing Road Speed, CIKM 2020 (KAIST)*

# Q&A