

# Preprocessing Pipeline

Ki Hyun Kim

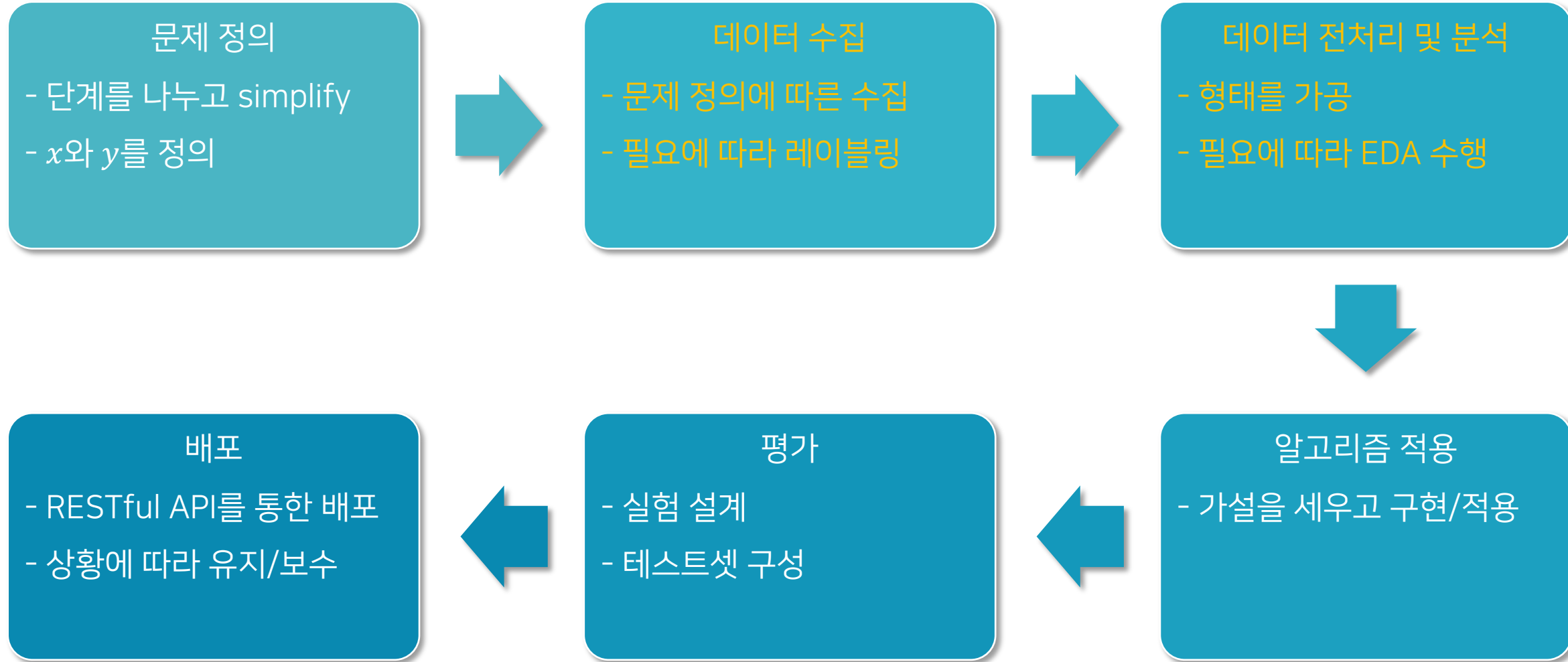
[nlp.with.deep.learning@gmail.com](mailto:nlp.with.deep.learning@gmail.com)

# 전처리의 늪에 오신 것을 환영합니다.

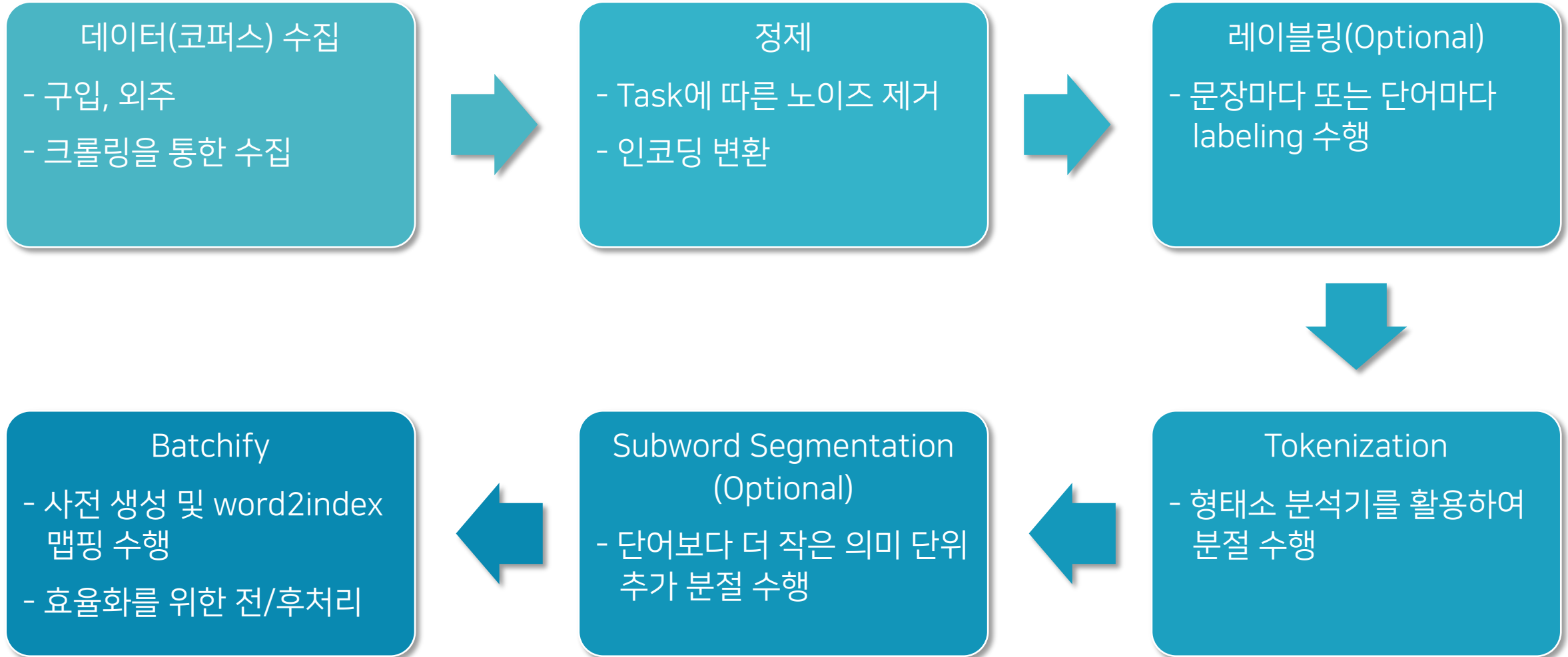
- 가장 재미없고(?) 반복적인 끝이 없는 작업
- 하지만 가장 중요 - 어쩌면 모델링만큼



# NLP Project Workflow



# Preprocessing Workflow



# 말뭉치(Corpus)란?

- 자연어처리를 위한 문장들로 구성된 데이터셋
- 복수표현: Corpora
- 포함된 언어 숫자에 따라,
  - Monolingual Corpus
  - Bi-lingual Corpus
  - Multilingual Corpus
- Parallel Corpus: 대응되는 문장 쌍이 labeling 되어 있는 형태

| English                 | Korean             |
|-------------------------|--------------------|
| I love to go to school. | 나는 학교에 가는 것을 좋아한다. |
| I am a doctor.          | 나는 의사 입니다.         |

# Service Pipeline

