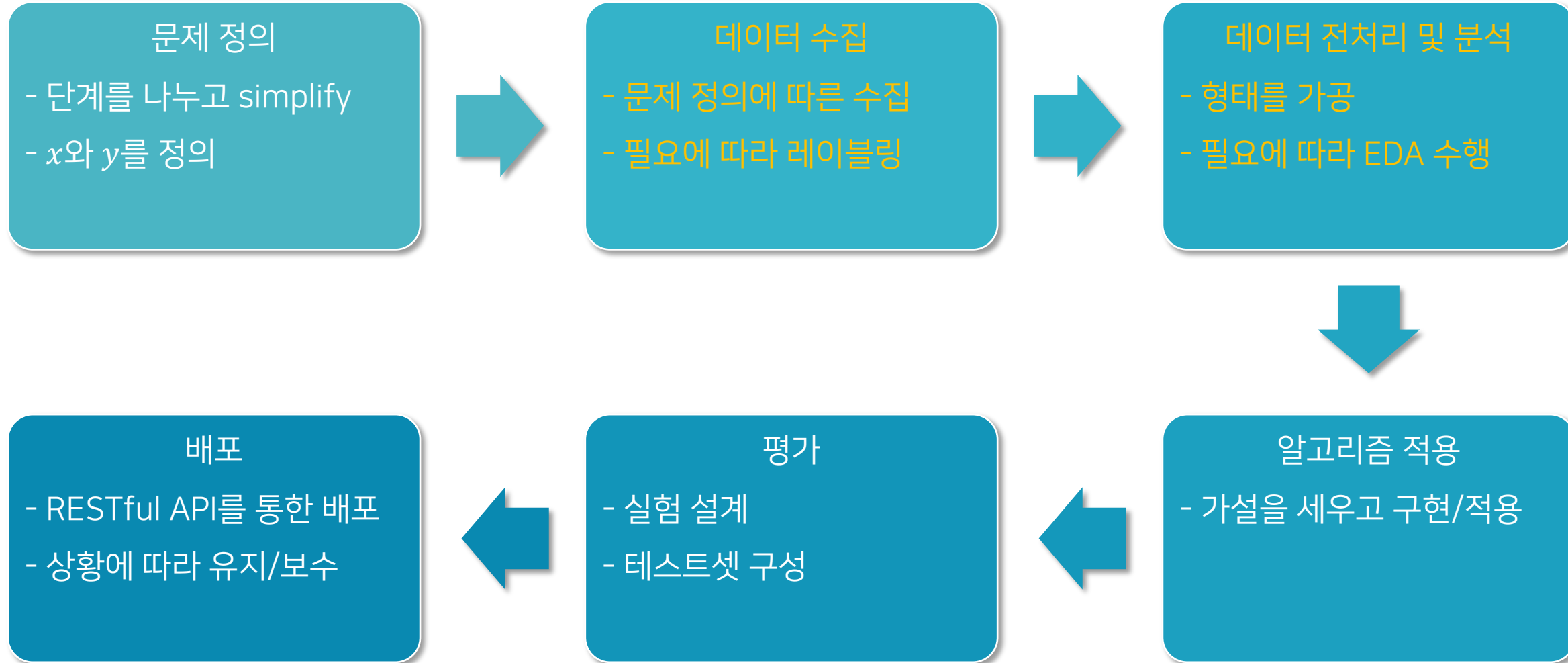


Preprocessing Summary

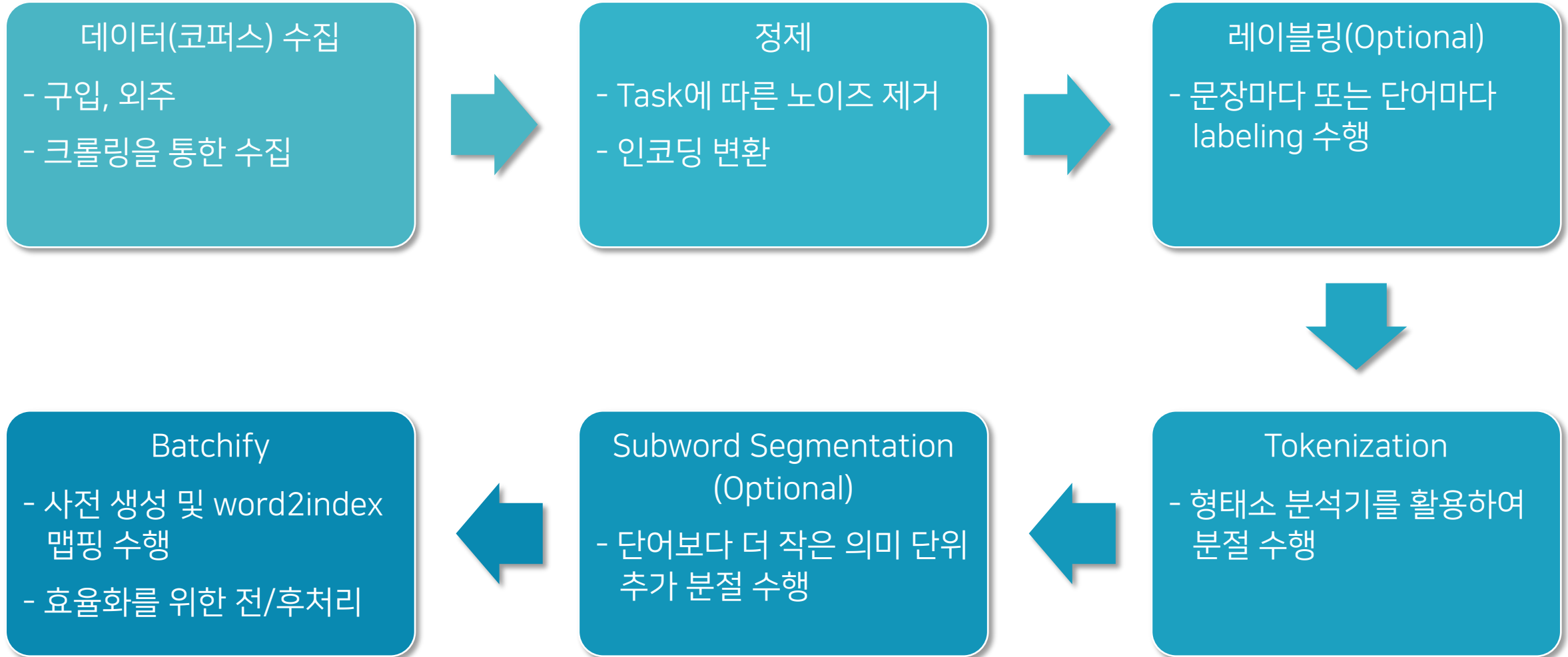
Ki Hyun Kim

nlp.with.deep.learning@gmail.com

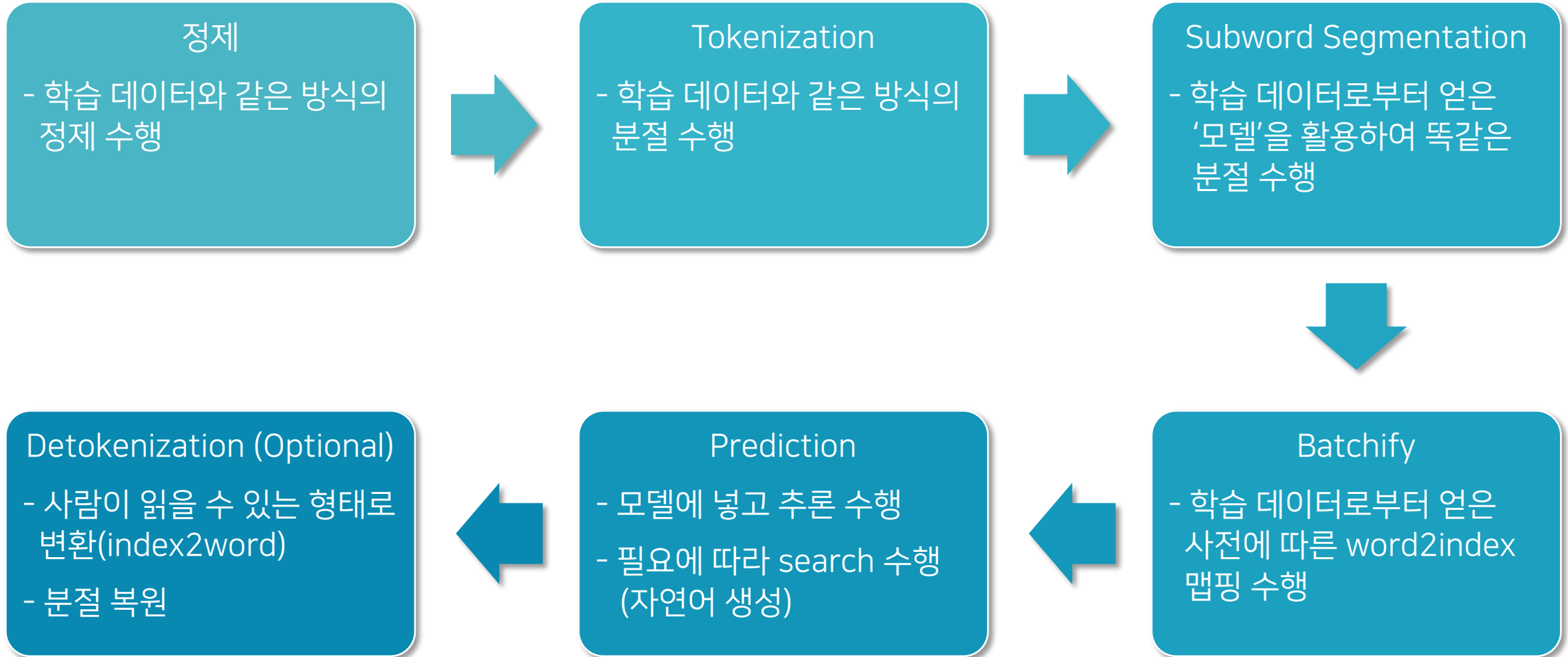
NLP Project Workflow



Preprocessing Workflow



Service Pipeline



Summary

- 정제
 - Task와 언어 및 도메인에 따른 특성
 - 풀고자 하는 문제의 특성에 따라 전처리 전략이 다름
 - 끝이 없는 과정
 - 노력과 품질 사이의 trade-off
 - Sweet spot을 찾아야함
- 분절
 - 한국어의 경우 띄어쓰기 normalization을 위해 형태소 분석기 활용이 필요
 - Subword segmentation을 통해 좀 더 잘게 분절 할 수 있음
- 모두 비슷한 알고리즘을 사용하고 있으므로, 결국 데이터의 양과 품질이 좌우함
 - 따라서 전처리 과정을 경시해서는 안됨