

챗봇 1만 개의 모델 서빙하기 : AI 서비스 어디까지 해봤니

고석현
Clova AI Business

NAVER

CONTENTS

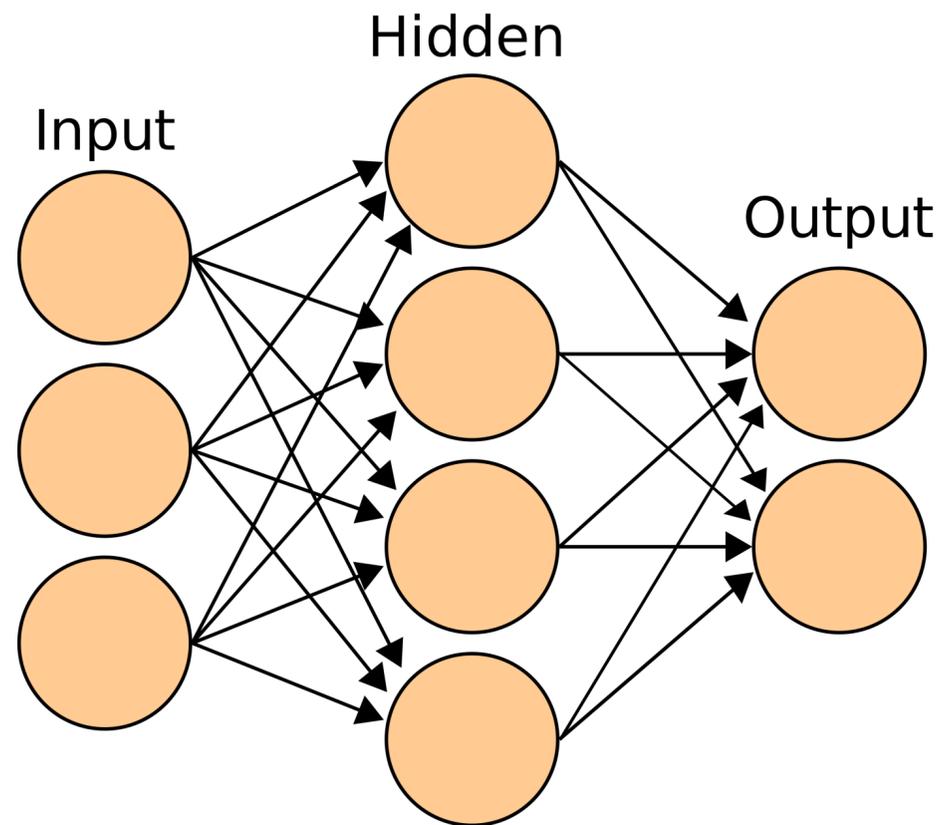
DEVIEW
2019

1. 모델의 성격 파악하기
2. Amdahl's law 를 기억하며
3. 조금만 더 줄어 주면 안될까?
4. 모델 정말 1만개를 서비스 했을까

1.모델의 성격 파악하기

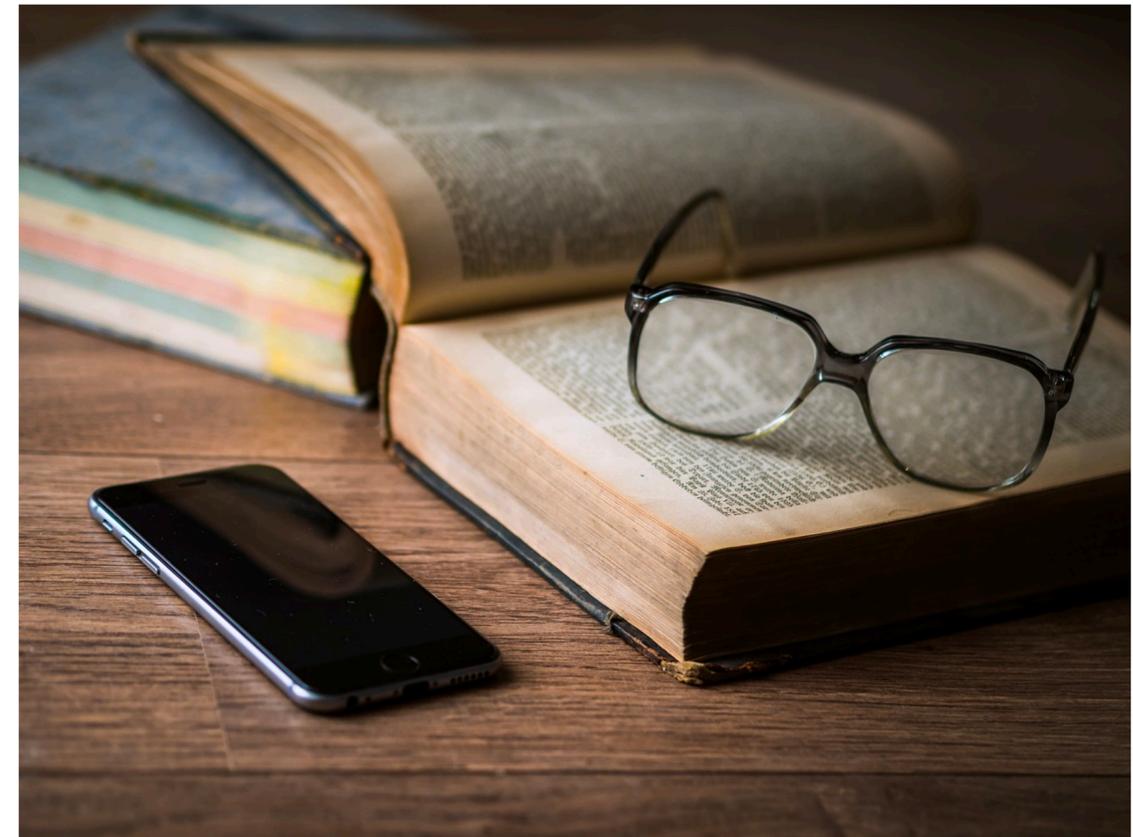
1.1 텍스트 ? 이미지 ? 음성 ?

실제로 중요한건 모델의 규모와 서비스 형태



1.1 텍스트 ? 이미지 ? 음성 ?

서비스 요구 사항에 따라 서버, 모바일, 임베디드등 여러 환경 고려 필요



1.2 안되는것을 빠르게 포기하는 방법

데이터 없이 학습 없이 이루어지는 마술같은 ML

1. 수백 ~ 수천 토큰의 질의를 이해하는 모델
2. 도메인 전용 학습, 튜닝이 필요하지 않는 범용 모델
3. 10턴 이상의 문맥을 고려하여 이해하는 모델

1.3 되는걸 빠르게 파악하기

구축 -> 훈련 -> 평가 -> 개선으로 이어지는 흐름을 자동화

1. 자동으로 데이터셋을 개선하기, 정리하기 (강재욱님 발표)
2. 자동으로 모델을 평가하고 튜닝하기 (2018 DEVIEW 이재원님 발표)
3. 분산 환경과 파이프 라인으로 최대한 빨리 모델 학습하기
4. 재학습, 운영, 관리, 지표화를 통해서 통제 가능하고 운영 가능한 머신러닝 서비스 만들기
5. 유저의 암시적 피드백을(IRF) 통해서 데이터 개선하기

1.4 모델의 비용 계산하기 1

- 베이스 라인과 모델 앙상블의 중요성
- 언어별, 데이터 규모별 강점을 가지는 다양한 모델이 존재
- SOTA 하나로 모든 데이터를 최적으로 서비스 하는것은 어려움



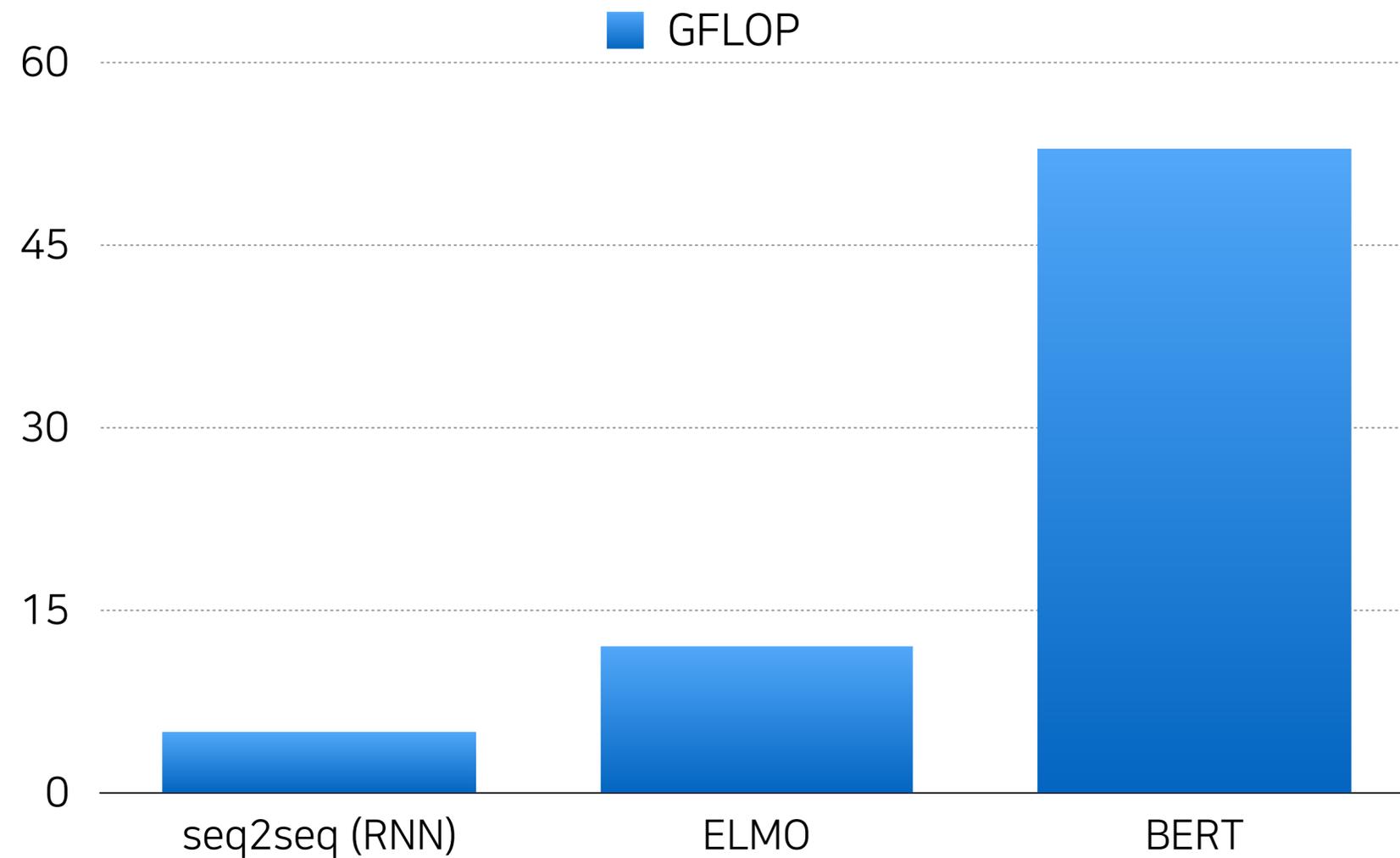
1.4 모델의 비용 계산하기 2

- 베이스 라인과 모델 앙상블의 중요성
- 언어별, 데이터 규모별 강점을 가지는 다양한 모델이 존재
- SOTA 하나로 모든 데이터를 최적으로 서비스 하는것은 어려움



1.4 모델의 비용 계산하기 3

- BERT(Transformer) 계열과 LSTM(RNN) 계열 모델의 계산 비용 비교



1.4 모델의 비용 계산하기 4

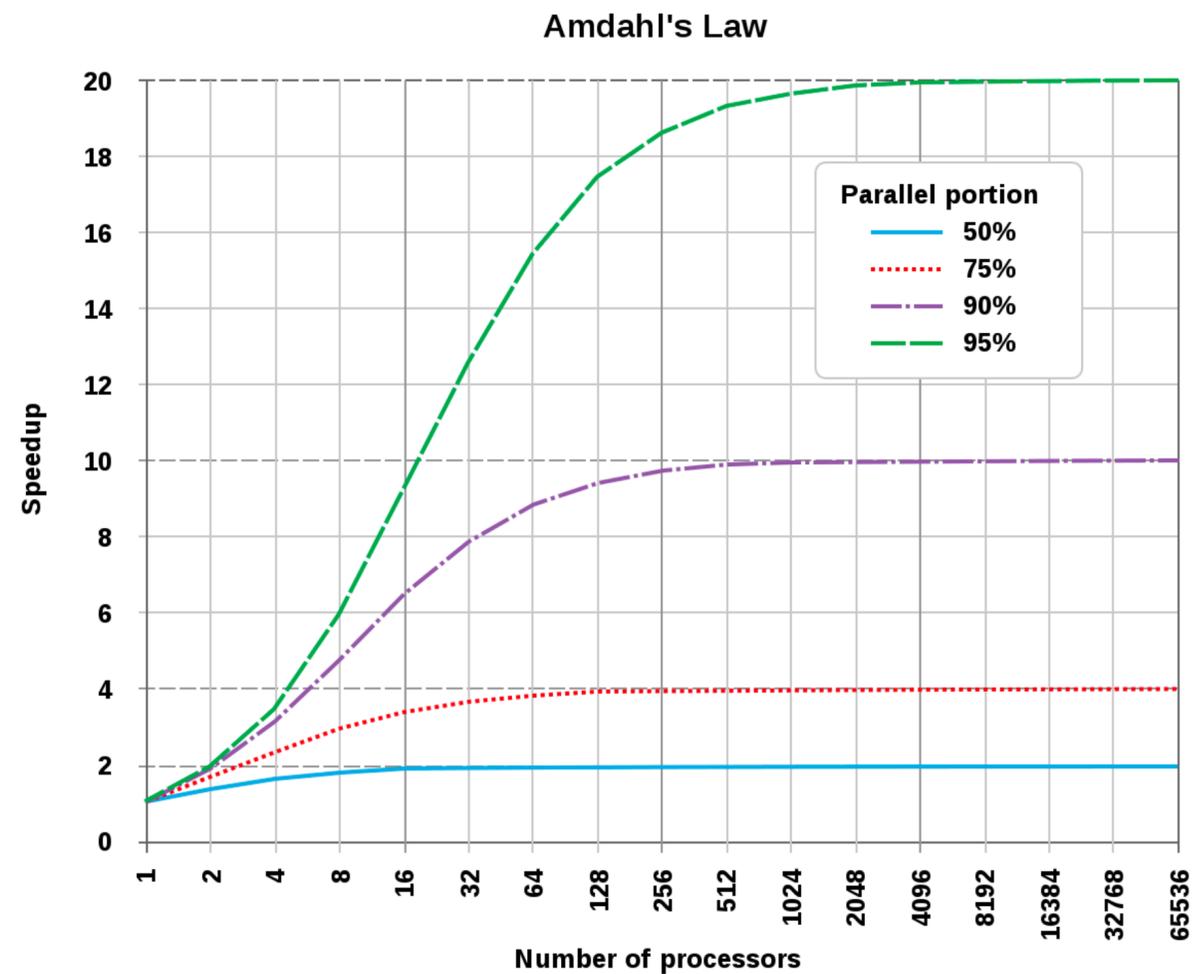
인공지능도 돈을 벌어야 합니다. 이제는 비용을 이야기 해야 합니다.

1. 모델 학습에 필요한 비용이 저렴하고
 2. 사전 학습된 모델을 도메인 별 fine-tuning 하여 이용 가능하며
 3. 모델의 연산 비용이 충분히 작고 최적화 되어 있으며
 4. 모델의 물리적인 용량이 작아 한정된 인프라에서 다수를 서비스 할 수 있으며
 5. 동시에 충분한 성능(정확도)를 보장해야 합니다.
0. 즉 비용을 고려한 최적화가 필요합니다.

2. Amdahl's law 를 기억하며

2.0 Amdahl's law

무한의 컴퓨팅 자원을 투입 하더라도



Definition [\[edit \]](#)

Amdahl's law can be formulated in the following way:

$$S_{\text{latency}}(s) = \frac{1}{(1-p) + \frac{p}{s}}$$

where

- S_{latency} is the theoretical speedup of the execution of the whole task;
- s is the speedup of the part of the task that benefits from improved system resources;
- p is the proportion of execution time that the part benefiting from improved resources originally occupied.

Furthermore,

$$\begin{cases} S_{\text{latency}}(s) \leq \frac{1}{1-p} \\ \lim_{s \rightarrow \infty} S_{\text{latency}}(s) = \frac{1}{1-p} \end{cases}$$

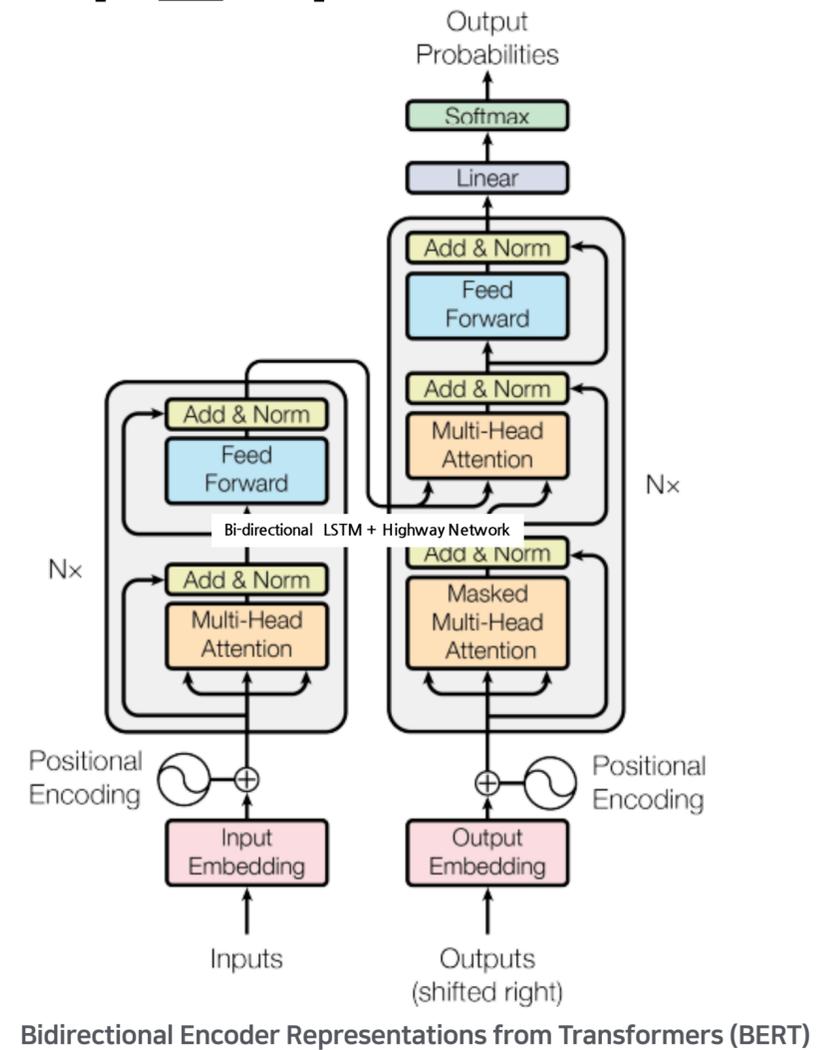
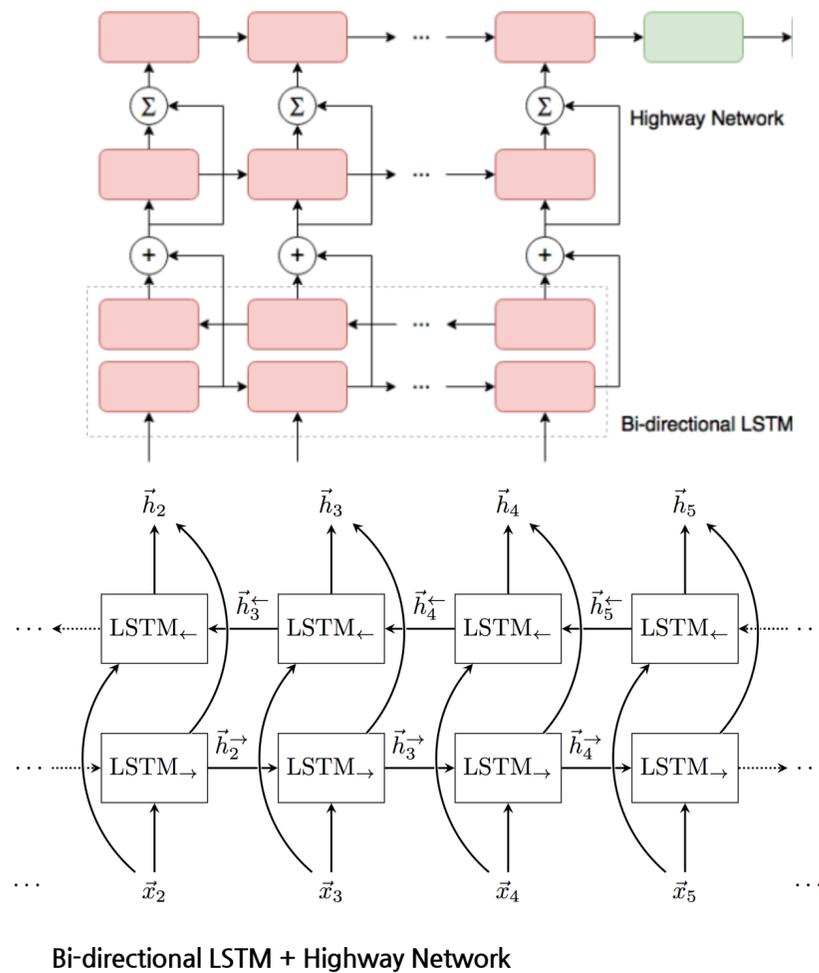
2.1 범용환경 ? 전용환경 ?

언제나 네이버 데이터 센터에 있다면 좋겠지만



2.2 빠른 모델이 좋은 모델일까 ?

계산 수행 시간과 총 계산 비용은 별도의 문제



2.3 아니 경제적인 모델

학습과 추론의 성능 최적화가 곧 비용인 시대

Artificial Intelligence / Machine Learning

Training a single AI model can emit as much carbon as five cars in their lifetimes

Deep learning has a terrible carbon footprint.

by Karen Hao

Jun 6, 2019

The artificial-intelligence industry is often compared to the oil industry: once mined and refined, data, like oil, can be a highly lucrative commodity. Now it seems the metaphor may extend even further. Like its fossil-fuel counterpart, the process of deep learning has an outsize environmental impact.

<https://www.technologyreview.com/s/613630/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>



3. 조금만 더 줄어 주면 안될까?

3.1 모델과 대화하기 그리고 설득하기

모델 레이어, 파라미터 튜닝

L1 가지 치기

레이어간 파라미터 공유

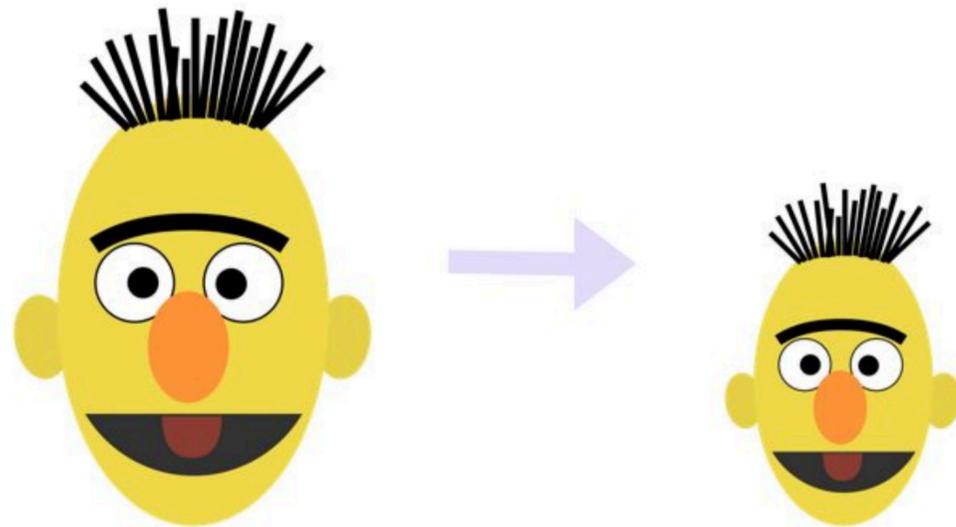
컴퓨팅 자원 친화적인 연산으로 교체

특정 장비별 오퍼레이션 최적화

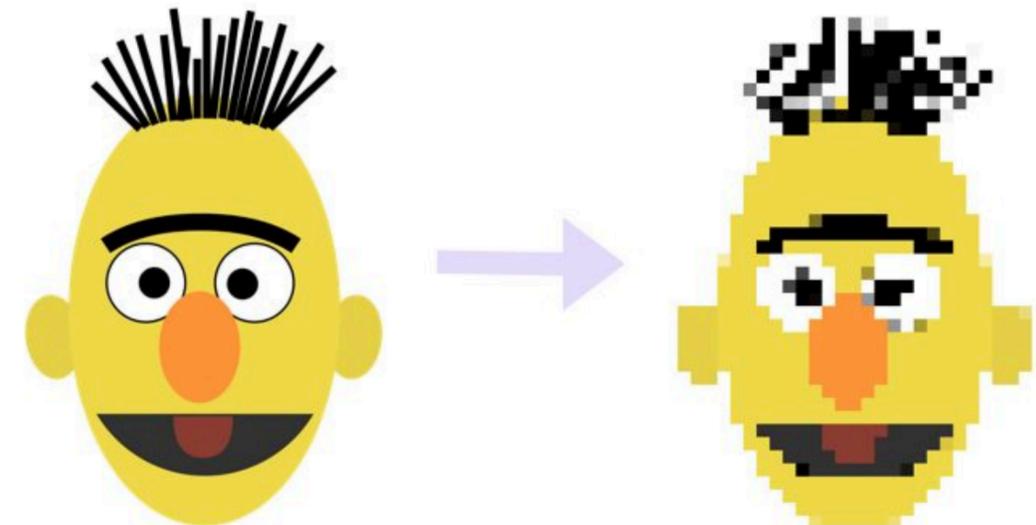
3.1 모델과 대화하기 그리고 설득하기

어떤 방법은 대가와 비용이 따릅니다.
모델의 정확도(성능)를 일부 포기해야 합니다.

Knowledge distillation

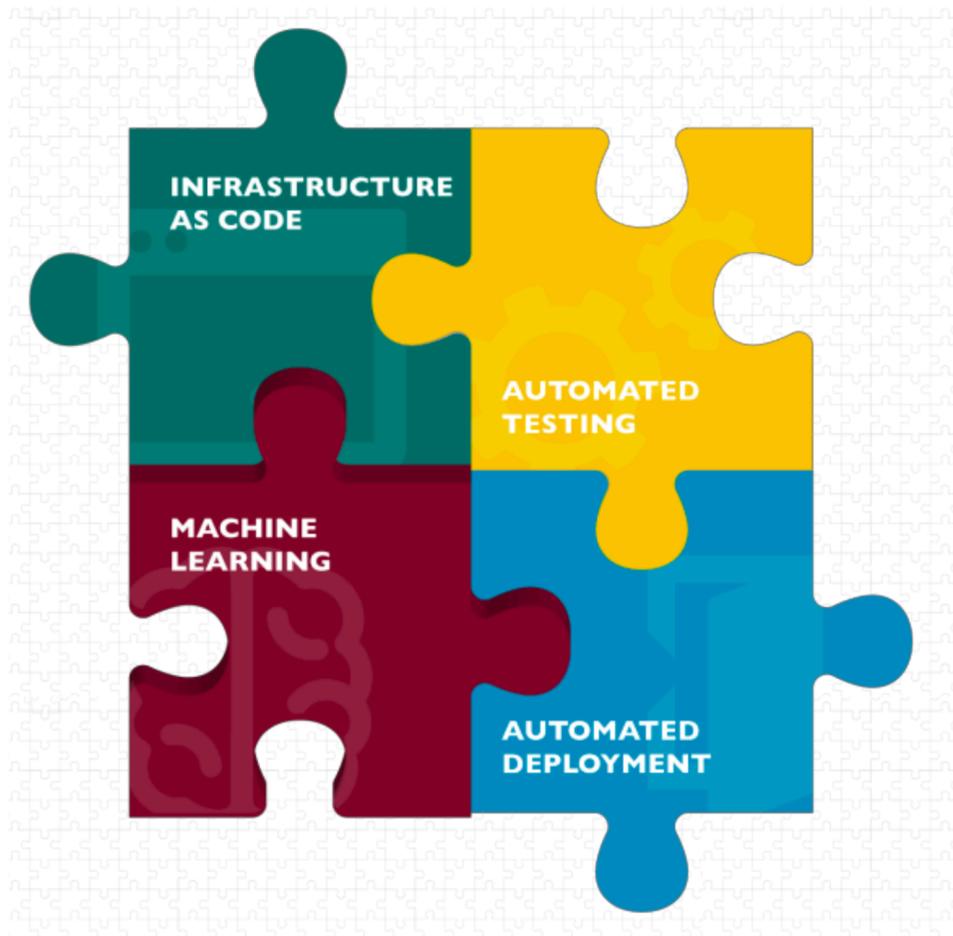


Quantization



3.2 그래도 안 될 거예요. 모델러와 소통하기

공동의 목표와 성과로 이어져야 하고 서비스와 제품으로 묶여야 합니다.



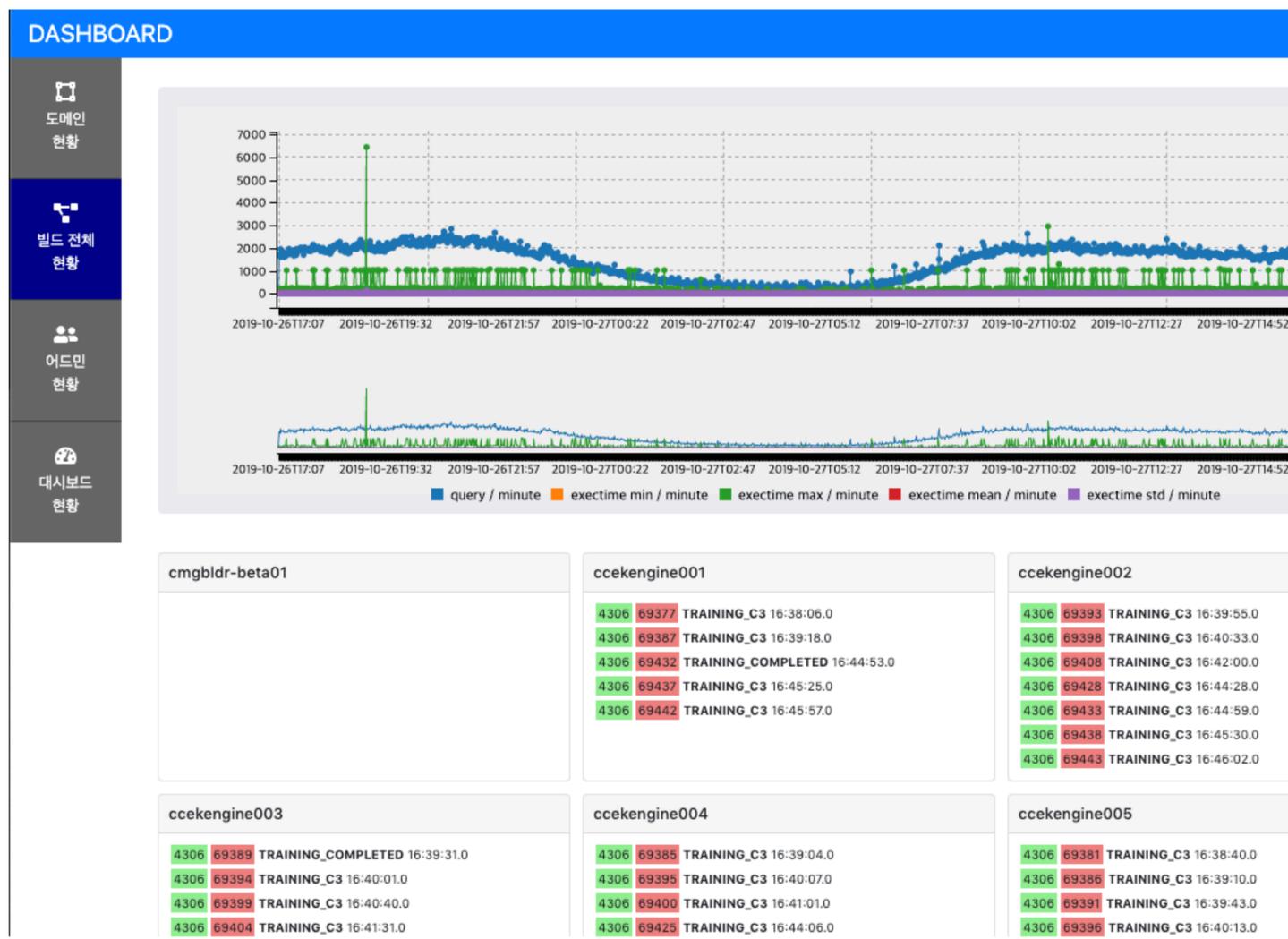
챗봇 빌더



3.2 그래도 안 될 거예요. 모델러와 소통하기

DEVIEW
2019

공동의 목표와 성과로 이어져야 하고 서비스와 제품으로 묶여야 합니다.



DASHBOARD KOR

real 4322 com.aiseminar.kevin CEK 모글리에서 보기 NCP 모글리에서 보기

신규 빌드 작업

현재 데이터를 기준으로 새로운 빌드를 추가할 수 있습니다. 버튼을 누르면 빌드가 됩니다. 호스트를 지정하면 선택된 해당 서버에서 빌드가 됩니다. 수정 배포 기능도 여기서 선택된 호스트의 영향을 받습니다.

호스트 지정 (기본) 빌드

빌드 현황

Copy CSV Excel PDF Print Show 10 entries Search:

빌드 ID	도메인 ID	생성일	C3DL ID	상태	DT	리얼 배포	베타 배포	개발 배포	개발 상태
63025	4322	2019-10-17 22:58:24.0	1570432872011_23564 1570432872011_23565	COMPLETED	2019-10-17T13-58-31.798Z	다시배포	다시배포	배포하기	수정배포
63024	4322	2019-10-17 13:55:23.0	1570432872011_23557 1570432872011_23558	COMPLETED	2019-10-17T13-55-29.819Z	배포하기	배포하기	배포하기	수정배포
62695	4322	2019-10-17 19:50:04.0	1570432872011_22995 1570432872011_23000	COMPLETED	2019-10-17T10-50-14.332Z	배포하기	배포하기	다시배포	수정배포
62659	4322	2019-10-17 17:56:11.0	1570432872011_22761 1570432872011_22764	COMPLETED	2019-10-17T08-56-17.518Z	배포하기	배포하기	배포하기	수정배포
62652	4322	2019-10-17 16:54:34.0	1570432872011_22645 1570432872011_22647	COMPLETED	2019-10-17T07-54-41.410Z	배포하기	배포하기	배포하기	수정배포
62497	4322	2019-10-17 13:34:08.0	1570432872011_22212 1570432872011_22214	COMPLETED	2019-10-17T04-34-10.435Z	배포하기	배포하기	배포하기	수정배포

배포 현황

Copy CSV Excel PDF Print Show 10 entries Search:

배포 ID	도메인 ID	빌드 ID	생성일	타겟	상태	DT	삭제
17945	4322	63025	2019-10-22 10:53:42.0	REAL	COMPLETED	/user/companyai/data/db/host=10.116.81.201/schema=mogli2tool/domainId=4322/dt=2019-10-17T13-58-31.798Z	삭제

3.2 그래도 안 될 거예요. 모델러와 소통하기

DEVIEW
2019

공동의 목표와 성과로 이어져야 하고 서비스와 제품으로 묶여야 합니다.



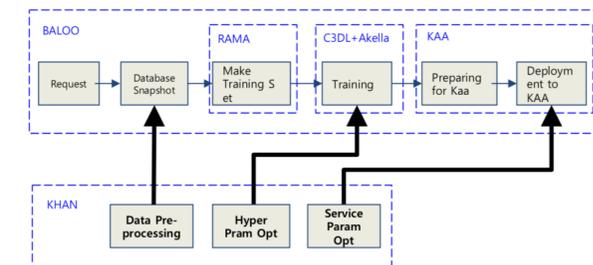
Khan

AutoML Project for Chatbot AI Builder

Overview

- End2End CB building pipeline
- Clean dataset
- Select appropriate features
- Select an appropriate models fam
- Optimize hyper-parameters
- Optimize service-parameters
- Automated model validation

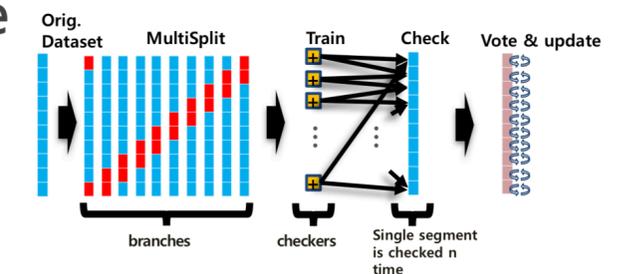
Chatbot Building Process



Automated Feature Eng.

- Intent and slot clustering
- Noisy query filtering
- Auto label error correction

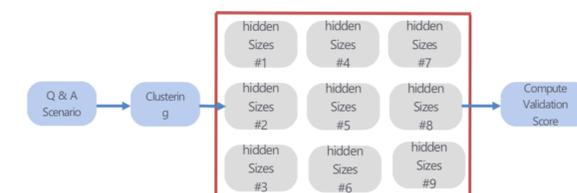
PICO: Auto Lable Correction



Hyper-Param Opt.

- Coordinate search
- Genetic algorithm
- Bayesian optimization

Coordinate Search for HiddenSizes

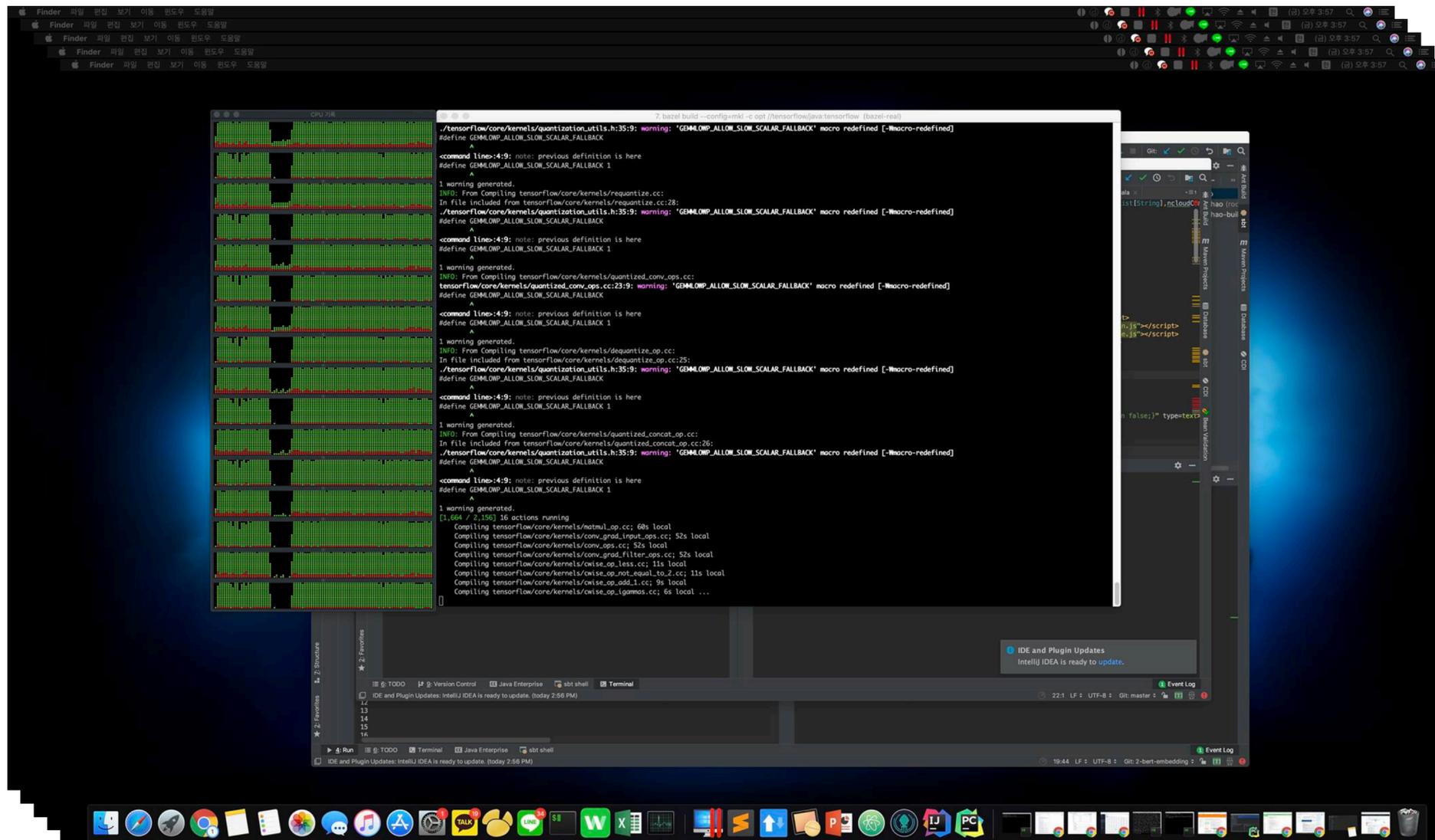


3.3 그래도 잘 안될거예요 이제 엔지니어에게 남은 모든걸 하기

DEVIEW
2019

프레임워크의 코드도 커스터마이징을 합니다.

수 없이 빌드와 테스트도 합니다. 연산 라이브러리도 교체합니다.

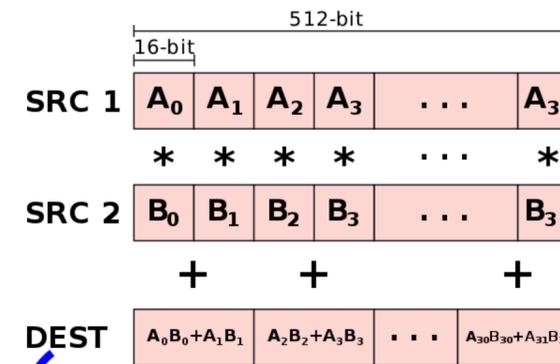


AVX-512 Vector Neural Network Instructions (VNNI) - x86

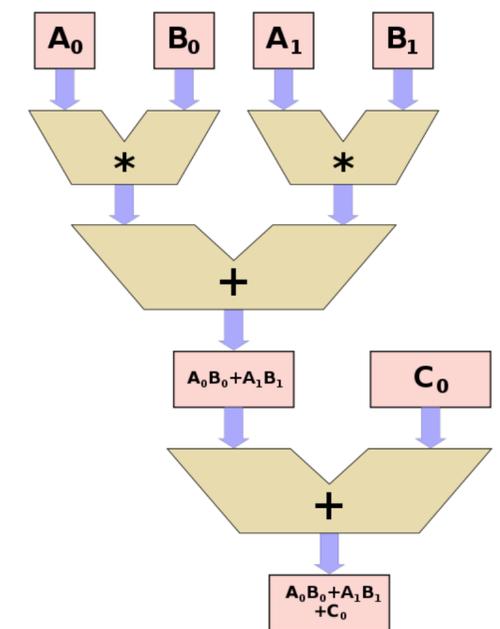
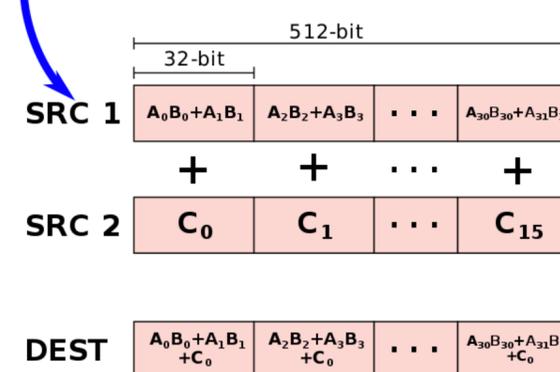
< x86

AVX-512 Vector Neural Network Instructions (AVX512 VNNI) is an x86 extension, algorithms.

VPMADDWD



VPADDD



3.3 그래도 잘 안될거예요 이제 엔지니어에게 남은 모든걸 하기

DEVIEW
2019

JNI(Java Native Interface) 수준으로 서비스 코드와 결합도 시도합니다.
네이티브 바이너리도 시스템과 환경별로 생성합니다.

```
import java.nio.IntBuffer
import org.tensorflow.{Graph, Session, Tensor}
import scala.collection.JavaConverters._

class TensorFlowProvider(model: Model) extends AutoCloseable {

  private val graph: Graph = {
    val graph = new Graph()
    graph.importGraphDef(model.getBytes)
    graph
  }

  private val session: Session = new Session(graph)
```

```
[redacted test tensorflow]$ ldd libtensorflow_framework.so.1
linux-vdso.so.1 => (0x00007fff0994c000)
libiomp5.so => /home/[redacted]/_bazel_irteam/e9fc3b58d0629b99
libmklml_intel.so => [redacted]e/bazel/_bazel_irteam/e9fc3b58d0
librt.so.1 => /lib64/librt.so.1 (0x00007fa2700b7000)
libpthread.so.0 => /lib64/libpthread.so.0 (0x00007fa26fe9b000)
libdl.so.2 => /lib64/libdl.so.2 (0x00007fa26fc97000)
libm.so.6 => /lib64/libm.so.6 (0x00007fa26f995000)
libstdc++.so.6 => /lib64/libstdc++.so.6 (0x00007fa26f68e000)
libgcc_s.so.1 => /lib64/libgcc_s.so.1 (0x00007fa26f478000)
libc.so.6 => /lib64/libc.so.6 (0x00007fa26f0ab000)
/lib64/ld-linux-x86-64.so.2 (0x00007fa27a947000)
[redacted test tensorflow]$ ldd java/libtensorflow_jni.so
linux-vdso.so.1 => (0x00007ffc27e47000)
libtensorflow_framework.so.1 (0x00007f0100c50000)
libiomp5.so => [redacted]/bazel/_bazel_irteam/e9fc3b58d0629b99
libmklml_intel.[redacted].cache/bazel/_bazel_irteam/e9fc3b58d0
libdl.so.2 => /lib64/libdl.so.2 (0x00007f00f8bc5000)
libm.so.6 => /lib64/libm.so.6 (0x00007f00f88c3000)
libpthread.so.0 => /lib64/libpthread.so.0 (0x00007f00f86a7000)
librt.so.1 => /lib64/librt.so.1 (0x00007f00f849f000)
libstdc++.so.6 => /lib64/libstdc++.so.6 (0x00007f00f8198000)
libgcc_s.so.1 => /lib64/libgcc_s.so.1 (0x00007f00f7f82000)
libc.so.6 => /lib64/libc.so.6 (0x00007f00f7bb5000)
/lib64/ld-linux-x86-64.so.2 (0x00007f010d2d0000)
```

3.3 그래도 잘 안될거예요 이제 엔지니어에게 남은 모든걸 하기

프로파일링을 통해 모델의 계산 비용 분석도 진행합니다.



<https://github.com/tensorflow/tensorflow/tree/master/tensorflow/core/profiler>

```
Run main
=====Model Analysis Report=====
Doc:
scope: The nodes in the model graph are organized by their names, which is hierarchical like filesystem.
flops: Number of float operations. Note: Please read the implementation for the math behind it.
Profile:
node name | # float_ops
_TFProfRoot | ██████████
tower0/ | /encoder/layer_5/output/dense/MatMul (150.99m/150.99m flops)
tower0/ | /encoder/layer_9/output/dense/MatMul (150.99m/150.99m flops)
tower0/ | /encoder/layer_5/intermediate/dense/MatMul (150.99m/150.99m flops)
tower0/ | /encoder/layer_0/output/dense/MatMul (150.99m/150.99m flops)
tower0/ | /encoder/layer_10/intermediate/dense/MatMul (150.99m/150.99m flops)
tower0/ | /encoder/layer_9/intermediate/dense/MatMul (150.99m/150.99m flops)
tower0/ | /encoder/layer_10/output/dense/MatMul (150.99m/150.99m flops)
```

3.3 그래도 잘 안될거예요 이제 엔지니어에게 남은 모든걸 하기

DEVIEW
2019

알려진 최적화들을 하나 하나씩 서비스에 적용합니다.

```
2019-10-15 22:00:16.824222: I tensorflow/core/platform/profile_utils/cpu_utils.cc:94] CPU Frequency: 2500000000 Hz
2019-10-15 22:00:16.828214: I tensorflow/compiler/xla/service/service.cc:168] XLA service 0x7f71cfbe9f20 initialized for platform Host (this does not guarantee that XLA will be used). Devices:
2019-10-15 22:00:16.828260: I tensorflow/compiler/xla/service/service.cc:176] StreamExecutor device (0): Host, Default Version
OMP: Info #212: KMP_AFFINITY: decoding x2APIC ids.
OMP: Info #210: KMP_AFFINITY: Affinity capable, using global cpuid leaf 11 info
OMP: Info #154: KMP_AFFINITY: Initial OS proc set respected: 0-79
OMP: Info #156: KMP_AFFINITY: 80 available OS procs
OMP: Info #157: KMP_AFFINITY: Uniform topology
OMP: Info #179: KMP_AFFINITY: 2 packages x 20 cores/pkg x 2 threads/core (40 total cores)
OMP: Info #214: KMP_AFFINITY: OS proc to physical thread map:
OMP: Info #171: KMP_AFFINITY: OS proc 0 maps to package 0 core 0 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 40 maps to package 0 core 0 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 4 maps to package 0 core 1 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 44 maps to package 0 core 1 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 8 maps to package 0 core 2 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 48 maps to package 0 core 2 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 65 maps to package 1 core 17 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 29 maps to package 1 core 18 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 69 maps to package 1 core 18 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 27 maps to package 1 core 19 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 67 maps to package 1 core 19 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 23 maps to package 1 core 20 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 63 maps to package 1 core 20 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 33 maps to package 1 core 24 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 73 maps to package 1 core 24 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 37 maps to package 1 core 25 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 77 maps to package 1 core 25 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 39 maps to package 1 core 26 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 79 maps to package 1 core 26 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 35 maps to package 1 core 27 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 75 maps to package 1 core 27 thread 1
OMP: Info #171: KMP_AFFINITY: OS proc 31 maps to package 1 core 28 thread 0
OMP: Info #171: KMP_AFFINITY: OS proc 71 maps to package 1 core 28 thread 1
OMP: Info #250: KMP_AFFINITY: pid 29837 tid 29947 thread 0 bound to OS proc set 0
2019-10-15 22:00:16.833505: I tensorflow/core/common_runtime/process_util.cc:115] Creating new thread pool with default inter op setting: 2. Tune using inter_op_parallelism_threads for best per
1
```

3.3 그래도 잘 안될거예요 이제 엔지니어에게 남은 모든걸 하기

DEVIEW
2019

시피유 데이터 시트를 뒤져봅니다. 물리적으로 연산장치 숫자가 적습니다. 인프라팀을 통해서 시제품 장비를 요청합니다.

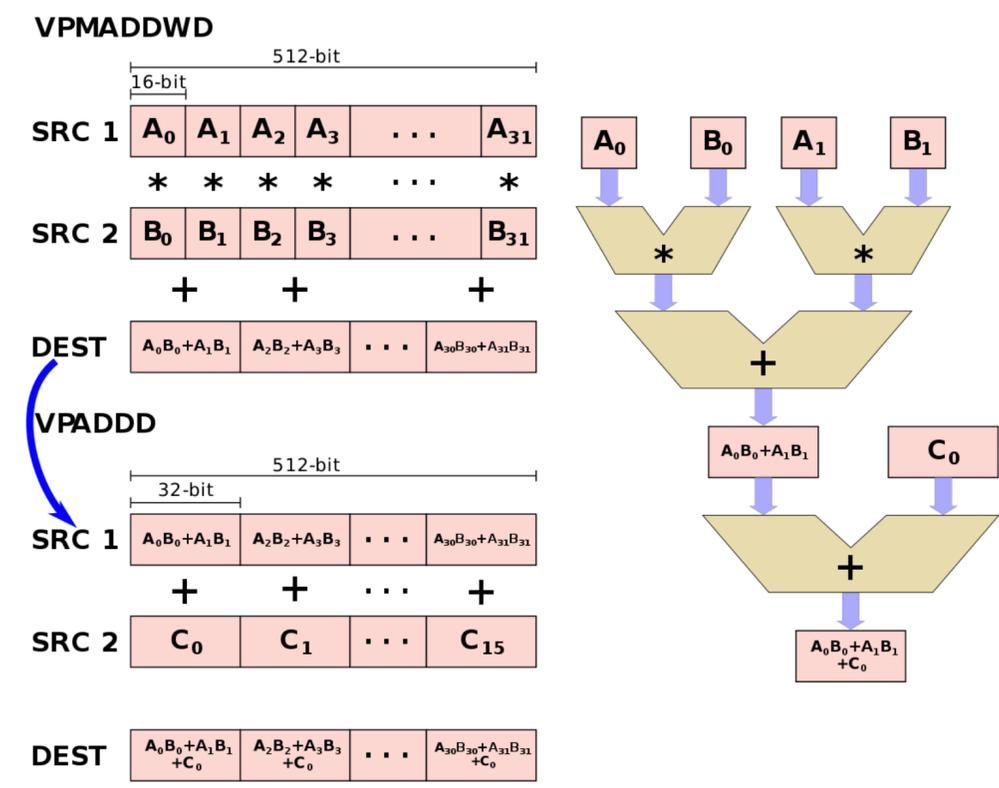
다시 테스트를 하고 아직 구현 최적화에 문제가 있는 부분도 찾습니다. intel과 미팅도 합니다.

Products Solutions Support			intel
Support Home > Product Specifications > Processors			
Advanced Technologies			
Intel® Optane™ Memory Supported ‡ ?	No	No	
Intel® Speed Shift Technology ?	Yes	Yes	
Intel® Turbo Boost Max Technology 3.0 ‡ ?	No	No	
Intel® Turbo Boost Technology ‡ ?	2.0	2.0	
Intel® vPro™ Platform Eligibility ‡ ?	Yes	Yes	
Intel® Hyper-Threading Technology ‡ ?	Yes	Yes	
Intel® Virtualization Technology (VT-x) ‡ ?	Yes	Yes	
Intel® Virtualization Technology for Directed I/O (VT-d) ‡ ?	Yes	Yes	
Intel® VT-x with Extended Page Tables (EPT) ‡ ?	Yes	Yes	
Intel® TSX-NI ?	Yes	Yes	
Intel® 64 ‡ ?	Yes	Yes	
Instruction Set Extensions ?	Intel® SSE4.2, Intel® AVX, Intel® AVX2, Intel® AVX-512 ?	Intel® SSE4.2, Intel® AVX, Intel® AVX2, Intel® AVX-512 ?	
# of AVX-512 FMA Units ?	1	2	
Enhanced Intel SpeedStep® Technology ?	Yes	Yes	
Intel® Volume Management Device (VMD) ?	Yes	Yes	

The screenshot shows a Git repository interface for tensorflow. The left sidebar displays the commit history with a list of commit messages, including updates to dockerfiles, fixes for missing activation methods, and improvements to the microfrontend. The main area shows a pull request with a list of changes, such as adding utility functions for SPIR-V operations and updating documentation. A specific commit is highlighted in the pull request list: '[Intel MKL] Upgrade curl to fix CVE-2019-5481 and CVE-2019-5482'.

3.4 AVX 512

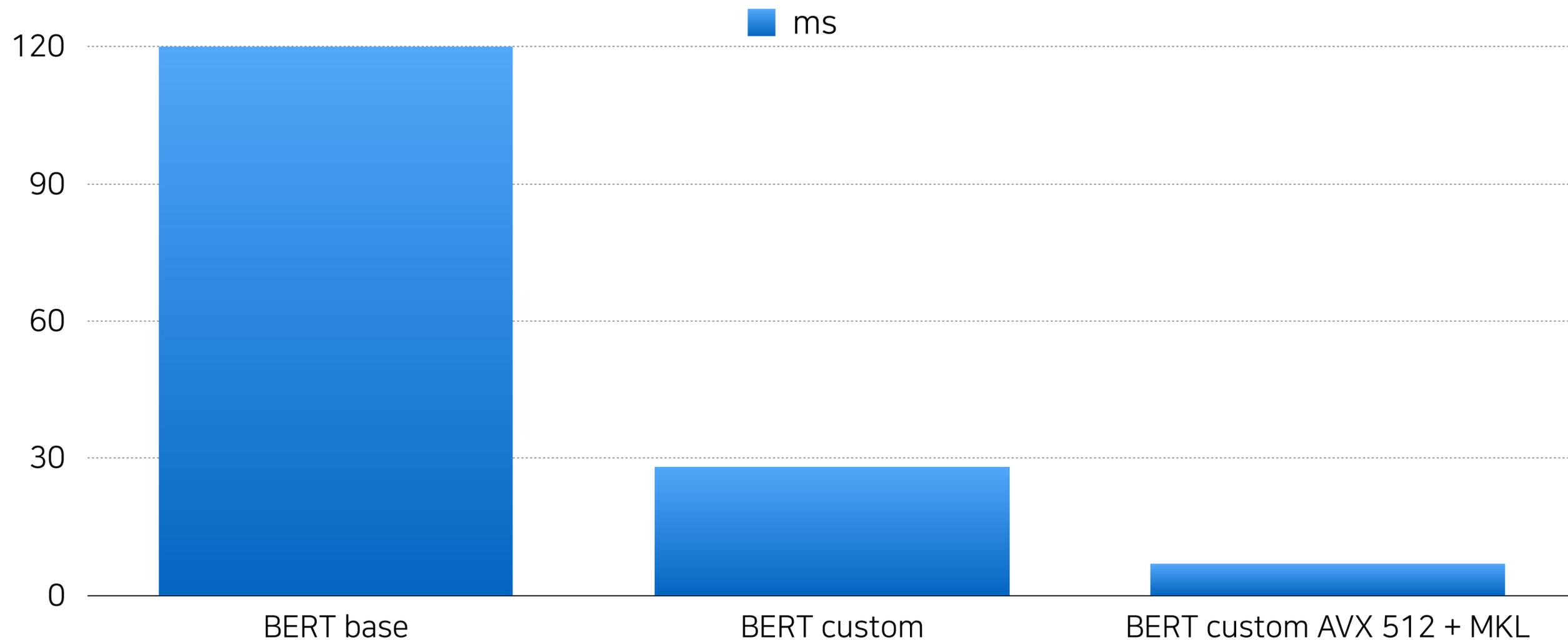
확장 레지스터와 512bit SIMD 명령어를 지원하는 AVX-512



3.4 AVX 512

현재 챗봇 서비스에 사용하는 모델 기준 단일 추론 7ms 수준으로 개선

- batch 1, intel xeon sp2 2.5Ghz 20 core



3.5 AVX 512 VNNI

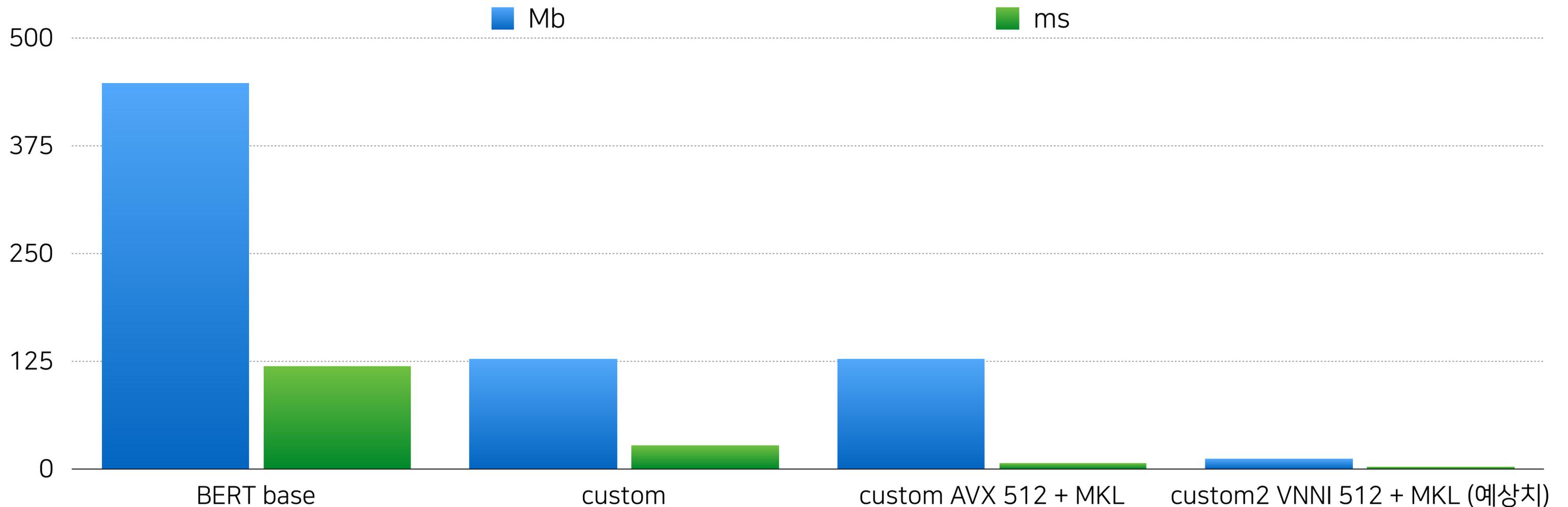
AVX-512 Vector Neural Network Instructions (VNNI)

- 전용 하드웨어의 개입 없는 없는 모델 양자화의 경우 실제 모델의 사이즈만 개선
- 모델 사이즈에 의한 메모리 제약 조건이 큰 모바일 탑재 환경에서만 유의미
- 하지만 AVX-512 VNNI등의 명령어는 하드웨어 수준에서 양자화된 모델 추론 가속

3.5 AVX 512 VNNI

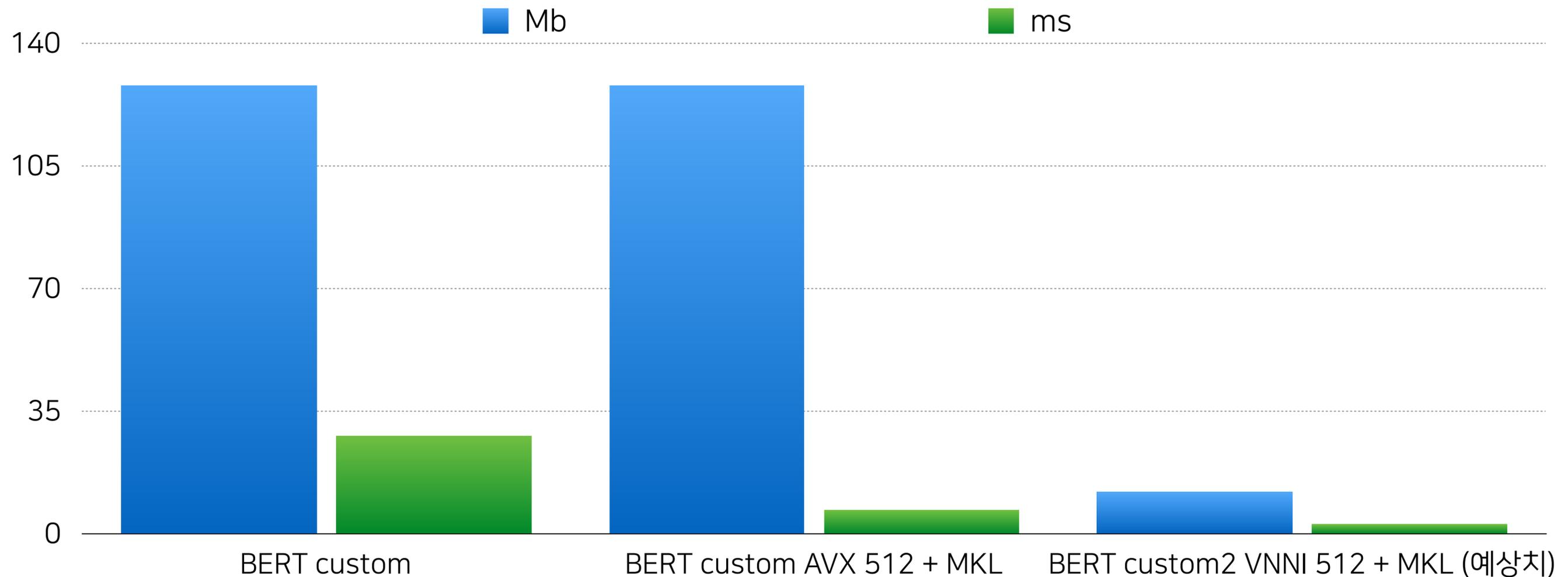
DEVIEW
2019

최신의 최적화 기법과 하드웨어 가속을 모두 적용한다면



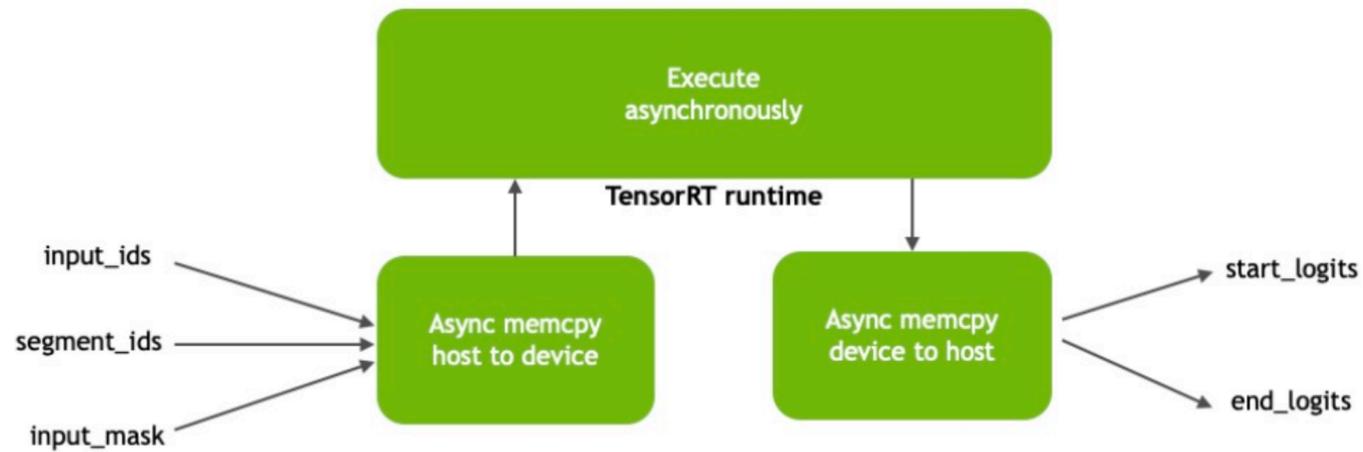
3.5 AVX 512 VNNI

최신의 최적화 기법과 하드웨어 가속을 모두 적용한다면

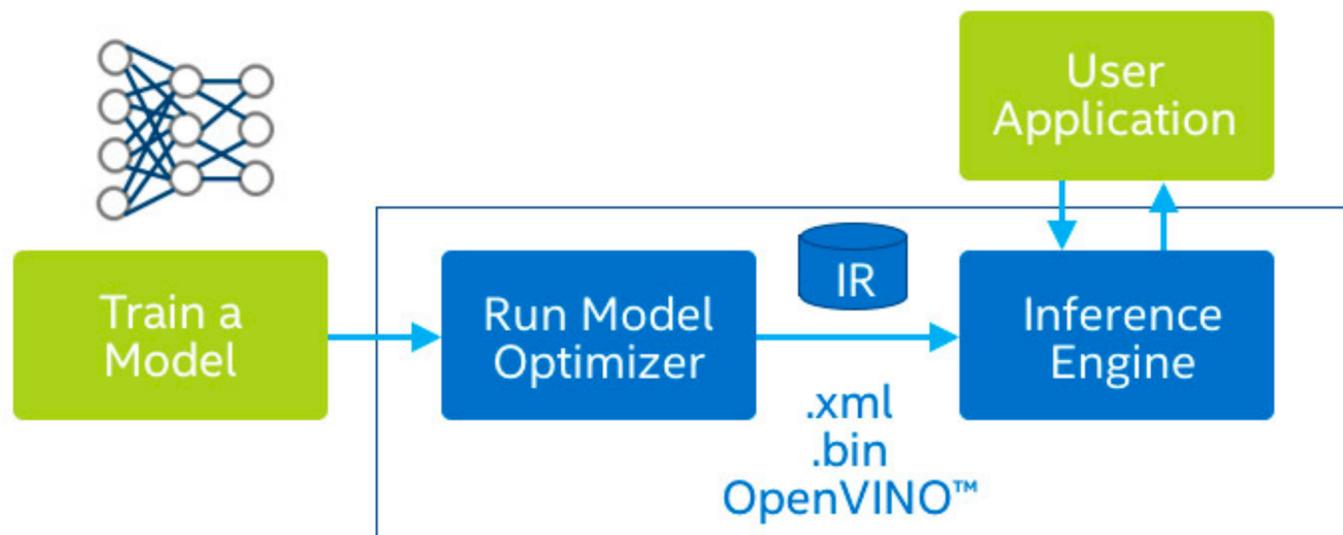


3.6 tensorRT

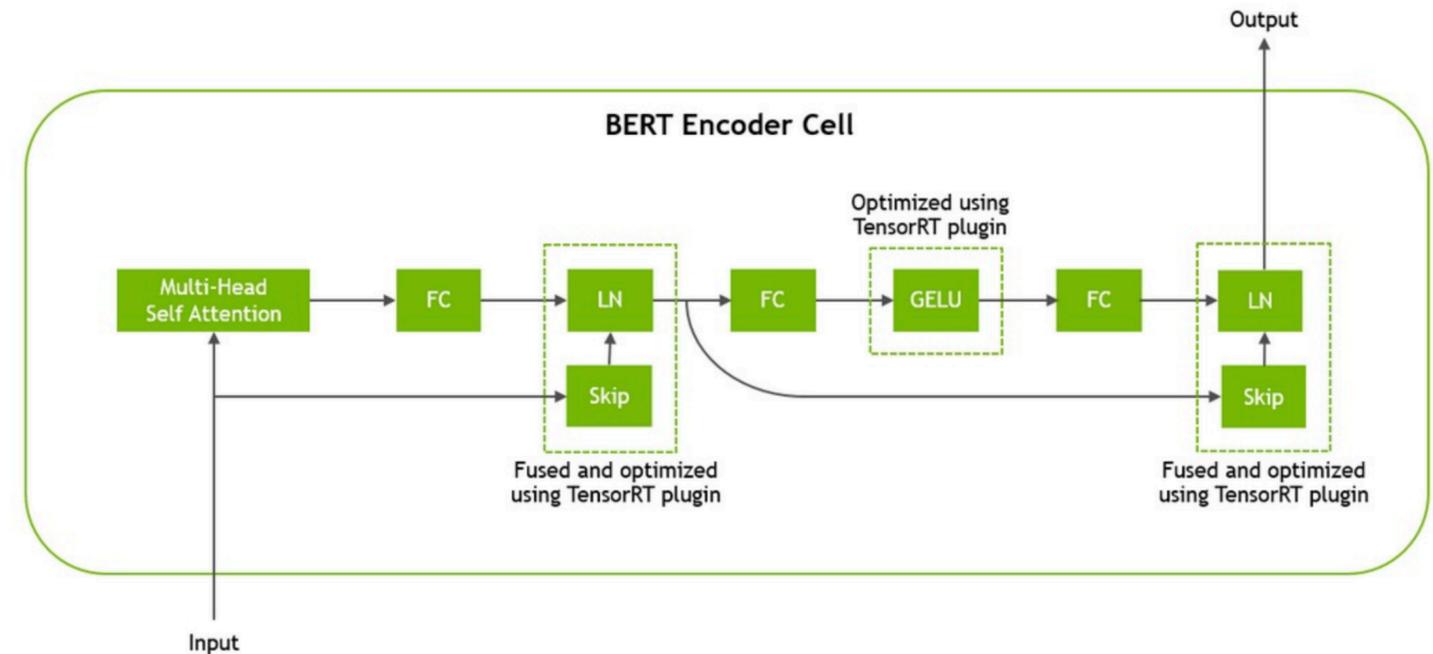
특정 환경 런타임을 기반으로한 추론 라이브러리등 존재



<https://developer.nvidia.com/tensorrt>



https://docs.openvino toolkit.org/latest/_docs_IE_DG_Introduction.html



4.모델 정말 1만개를 서비스 했을까

4.1 auto ml + auto quantization

Auto ML에서 사용 가능한 재료들을 최대한 자동화

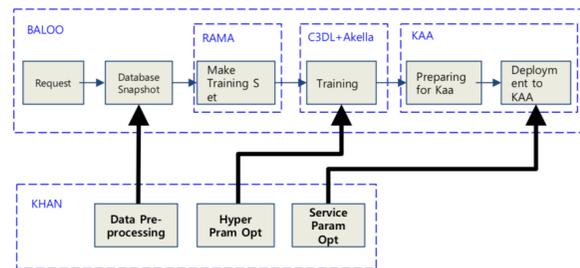
Khan

AutoML Project for Chatbot AI Builder

Overview

- End2End CB building pipeline
- Clean dataset
- Select appropriate features
- Select an appropriate models fam
- Optimize hyper-parameters
- Optimize service-parameters
- Automated model validation

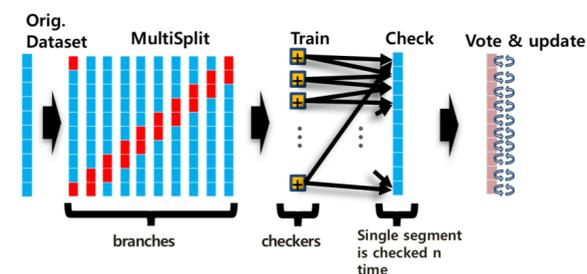
Chatbot Building Process



Automated Feature Eng.

- Intent and slot clustering
- Noisy query filtering
- Auto label error correction

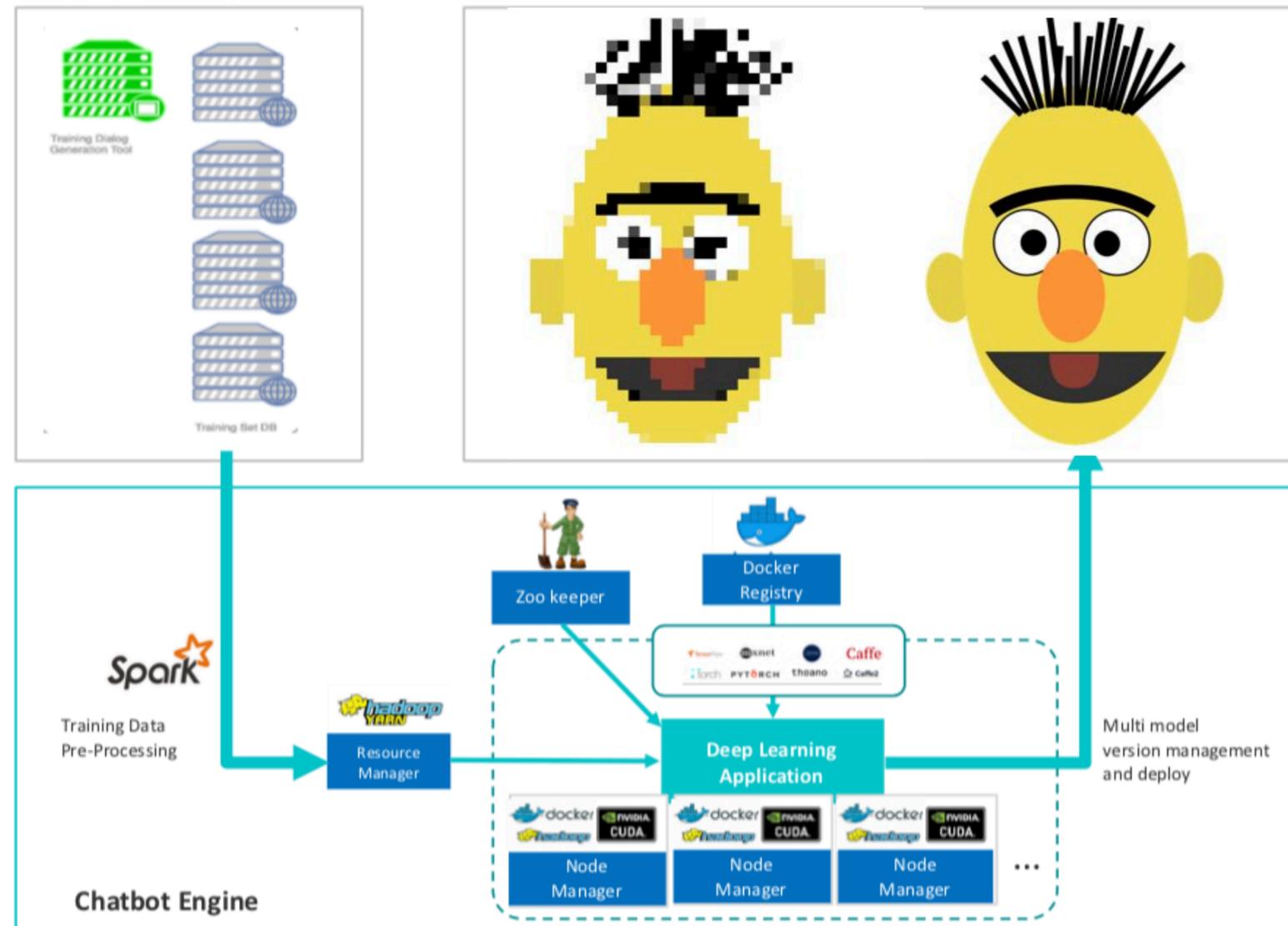
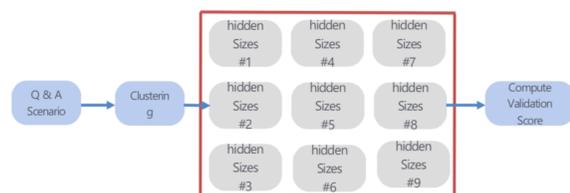
PICO: Auto Lable Correction



Hyper-Param Opt.

- Coordinate search
- Genetic algorithm
- Bayesian optimization

Coordinate Search for HiddenSizes



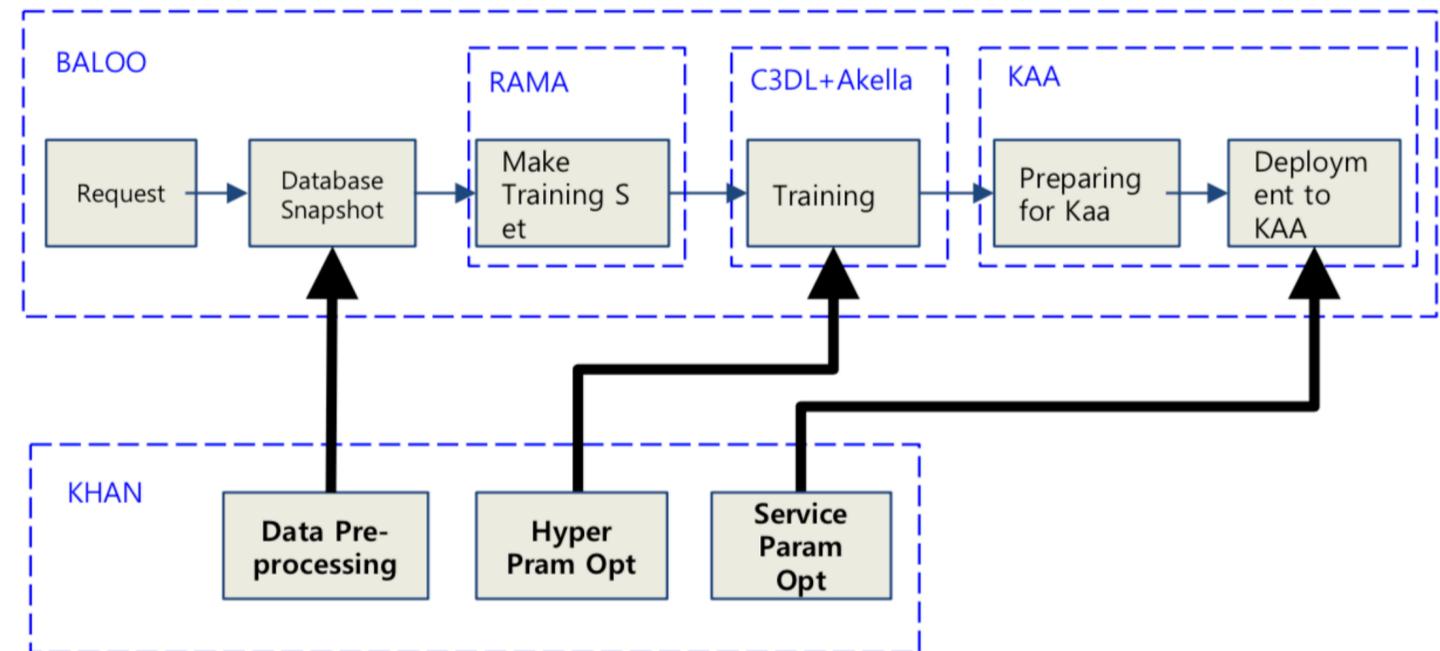
4.2 one source multi environment models

DEVIEW
2019

다양한 환경에서 (일반서버, GPU서버, 웨어러블, 가상, 오케스트레이션)
서빙 가능한 모델을 하나의 소스와 빌드 시스템으로 관리

ID	대화 이름	대화 유형	대화 위치	질문	질문 개수	답변	답변 유형	학습 상태	작성/변경일 (UTC+09:00)	수정
3167665	클로바 강의등록	클로바강...	단일 대화	네이버 챗봇 만들거요	8	네이버클로바에서 진행하는 챗봇 강의 증가...	기본 답변	학습 완료	2019-10-17 10:09:15	수정
3167664	강의변경	강의변경	단일 대화	다른 경로 변경해주세요	7	강의변경을 도와드리겠습니다. 변경을 원...	기본 답변	학습 완료	2019-10-17 19:35:55	수정
3167663	강의정보입력	강의정보...	단일 대화	챗봇만들거요	7	전화주신 휴대폰번호 끝 네 자리가 이상삼...	기본 답변	학습 완료	2019-10-17 17:54:30	수정
3158317	본인확인	공부정	단일 대화	네 맞아요	5	한시영 님, 확인 감사합니다. 아직 결제...	기본 답변	학습 완료	2019-10-17 17:53:43	수정
3158308	무료주차 가능여부	주차안내	단일 대화	혹시 무료주차 가능한가요?	21	신청자 본인에 한하여 일일 무료주차 ...	기본 답변	학습 완료	2019-10-17 10:09:58	수정

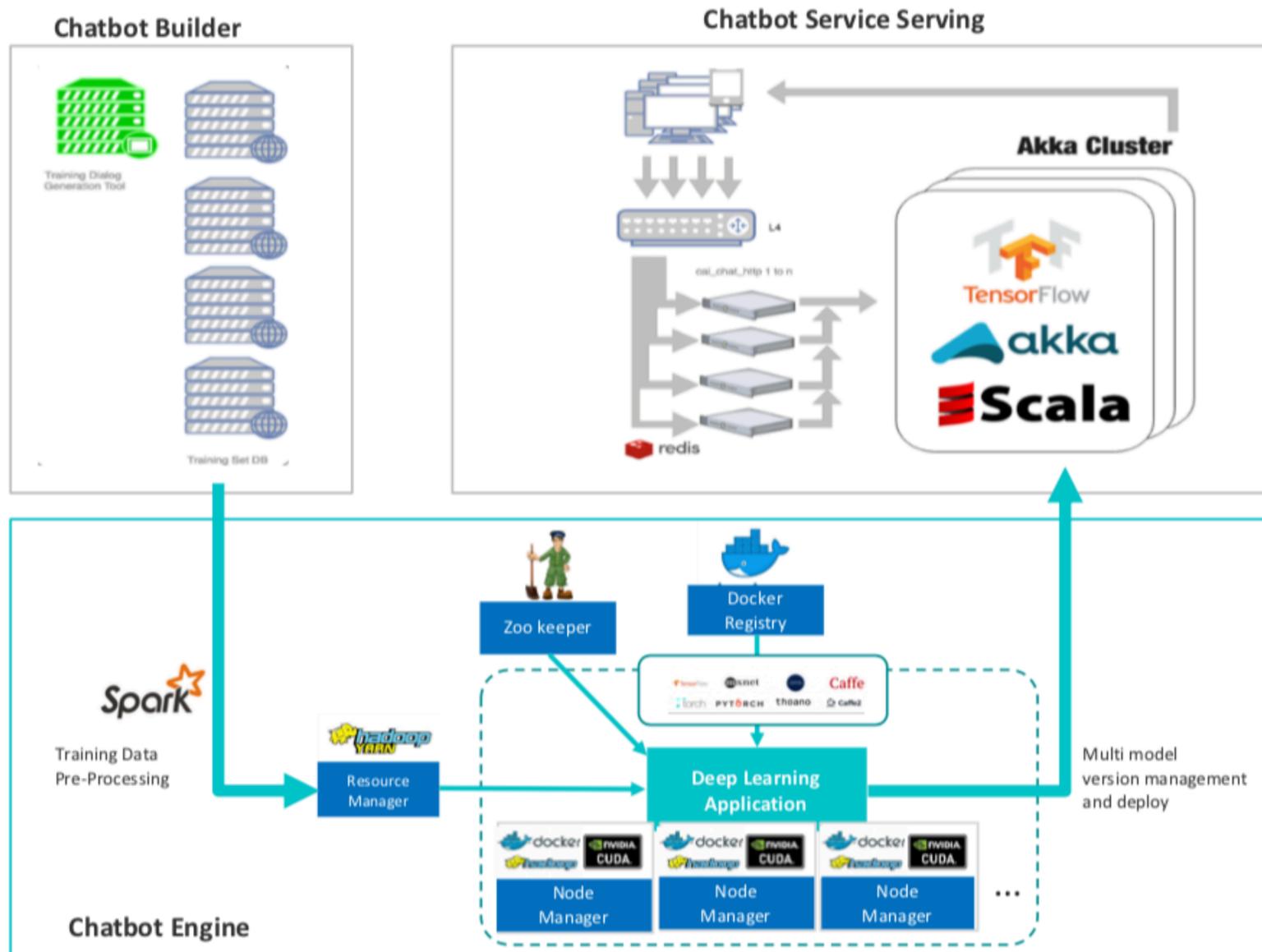
Chatbot Building Process



4.2 one source multi environment models

DEVIEW
2019

학습된 모델과 데이터를 프리즈된 버전별 스냅샷 통해 저장, 관리, 배포



신규 빌드 작업

현재 데이터를 기준으로 새로운 빌드를 추가할수 있습니다. 버튼을 누르면 빌드 됩니다.

호스트 지정 (기본)

빌드 현황

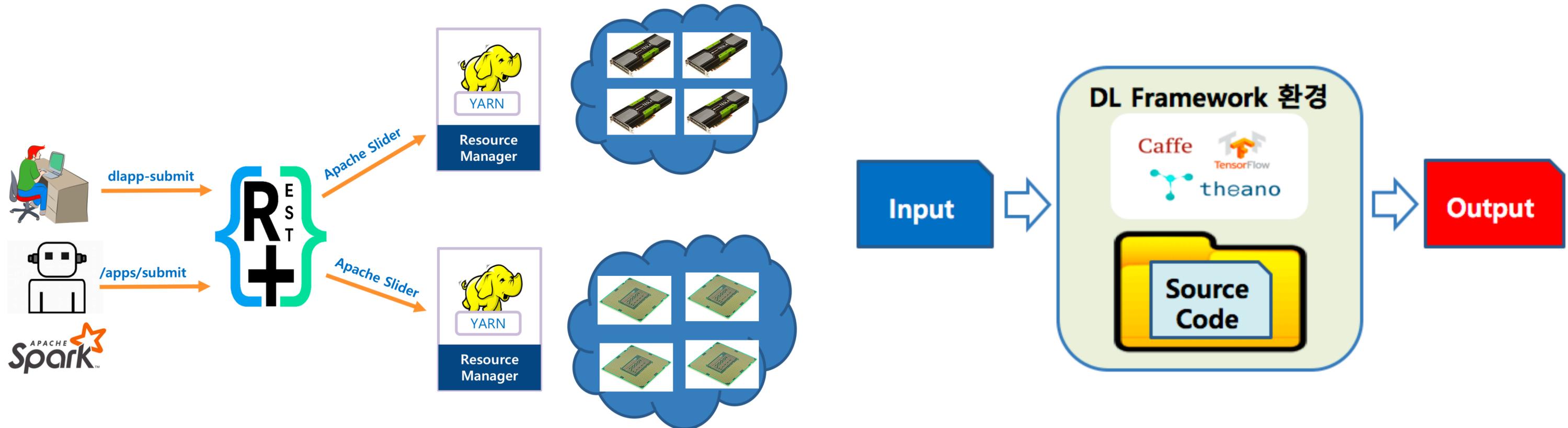
Copy CSV Excel PDF Print Show 10 entries

DT	리얼 배포	베타 배포	개발 배포	개발 상태
2019-10-17T13-58-31.798Z	다시배포	다시배포	배포하기	수정배포
2019-10-17T13-55-29.819Z	배포하기	배포하기	배포하기	수정배포
2019-10-17T10-50-14.332Z	배포하기	배포하기	다시배포	수정배포
2019-10-17T08-56-17.518Z	배포하기	배포하기	배포하기	수정배포
2019-10-17T07-54-41.410Z	배포하기	배포하기	배포하기	수정배포
2019-10-17T04-34-10.435Z	배포하기	배포하기	배포하기	수정배포

4.2 one source multi environment models

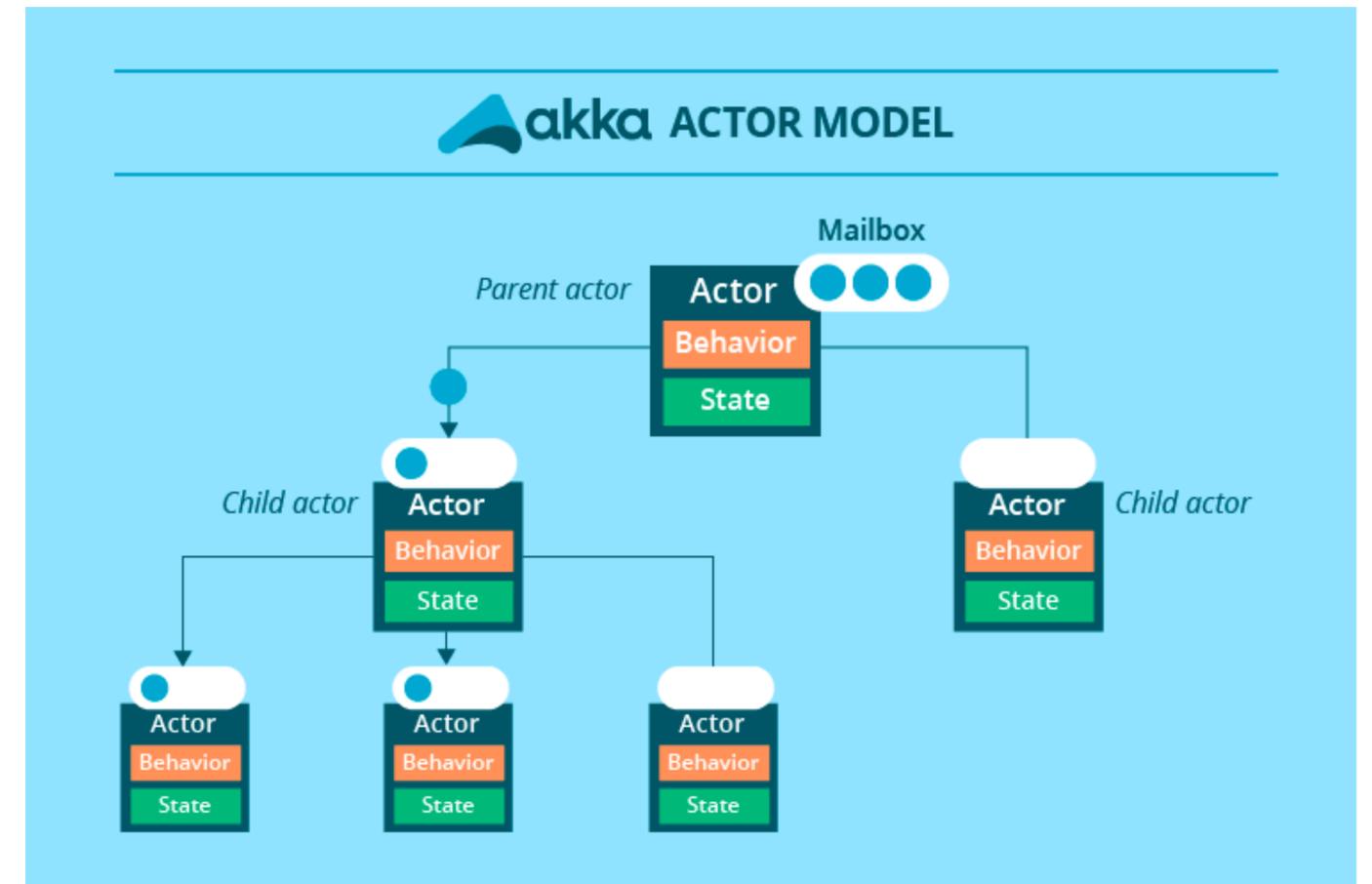
DEVIEW
2019

네이버 딥러닝 분산 플랫폼, C3DL 기반으로 동작



4.3 decentralized clusters - (kaa)

비 중앙화된 분산 클러스터를 이용한 머신러닝 모델 서빙 엔진을 설계



4.4 decentralized clusters - (kaa)

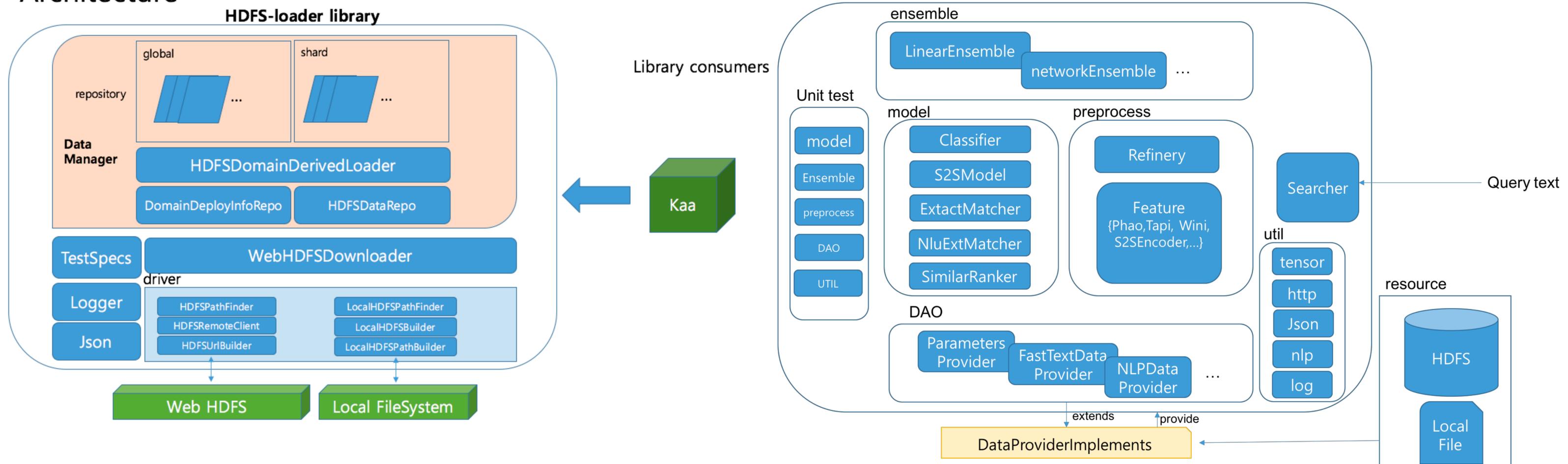
하나의 코드 베이스로 동작하는 데이터 파이프라이닝과 서빙 시스템



4.4 decentralized clusters - (kaa)

모델의 실험과 평가의 대한 공통 코드를 정적 모듈화 -> 신속한 모델 설계, 실험

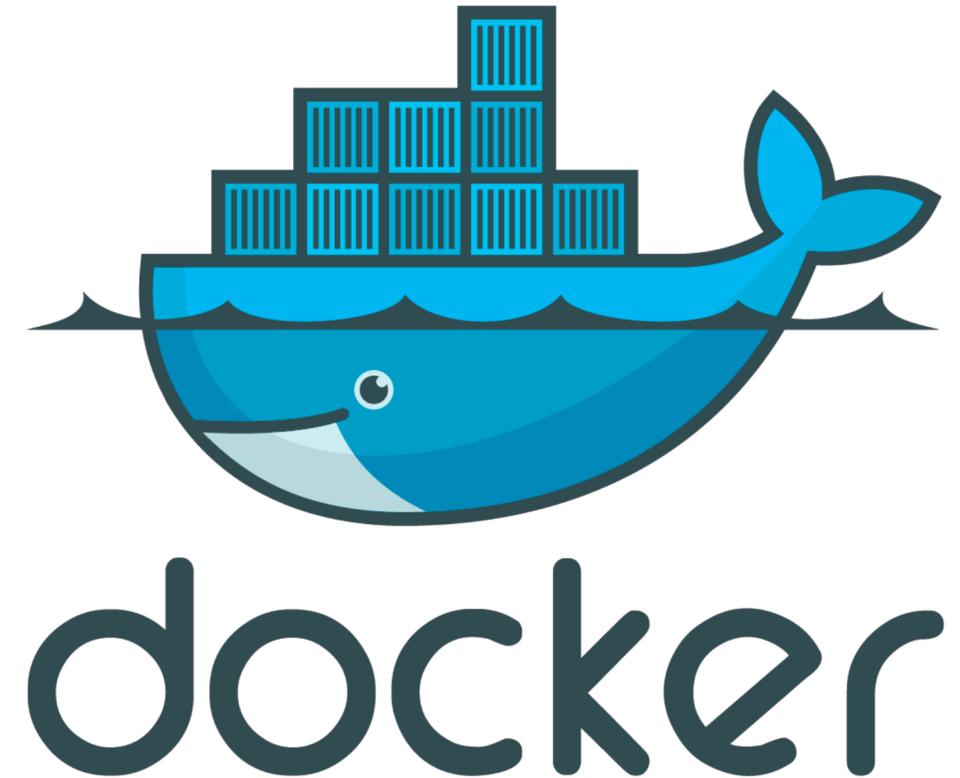
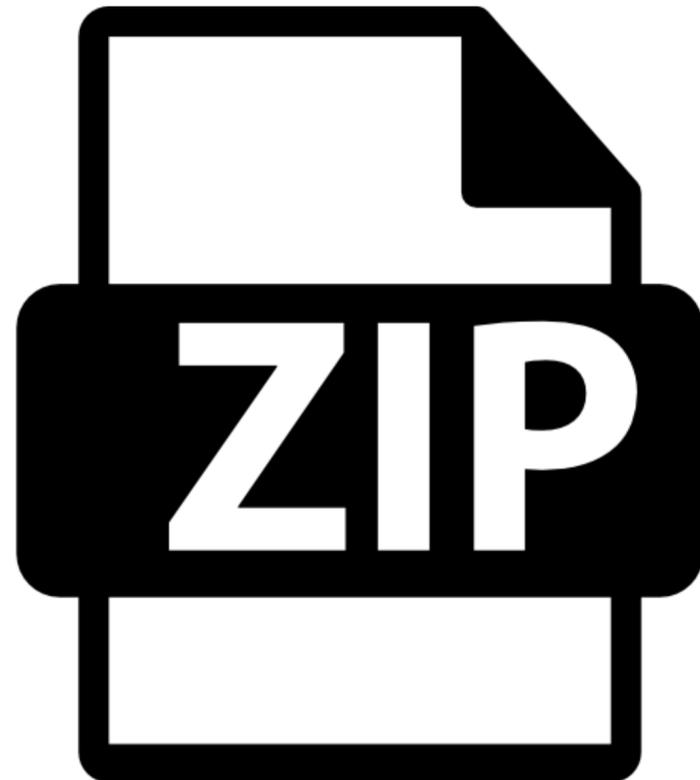
Architecture



4.5 decentralized clusters - (kaa)

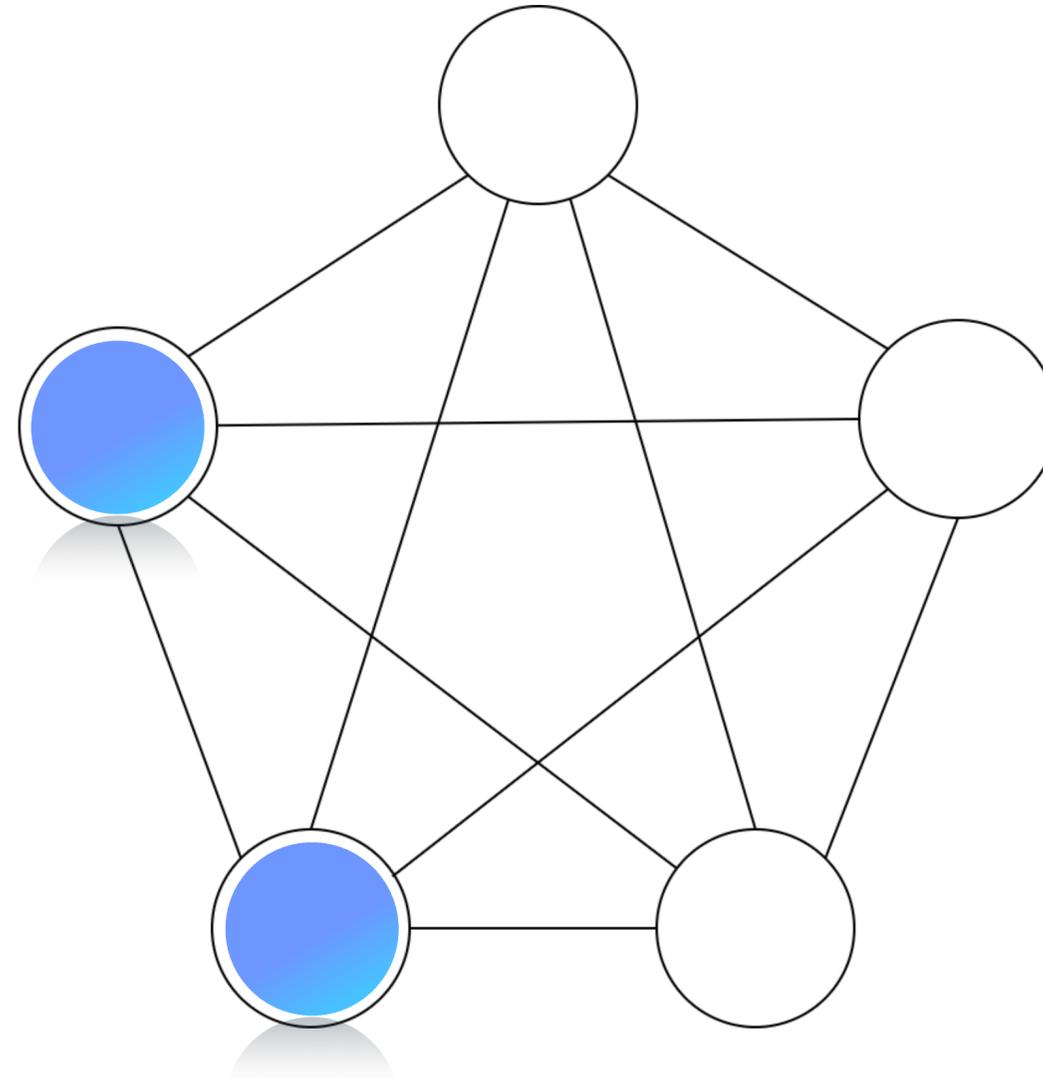
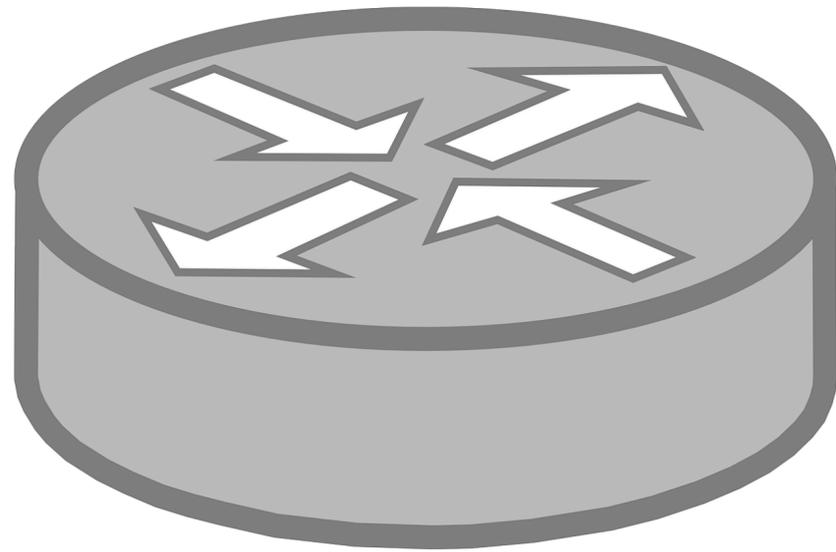
DEVIEW
2019

On-premise등 B2B를 위한 패키징을 시스템 초기 설계 수준에서 고려



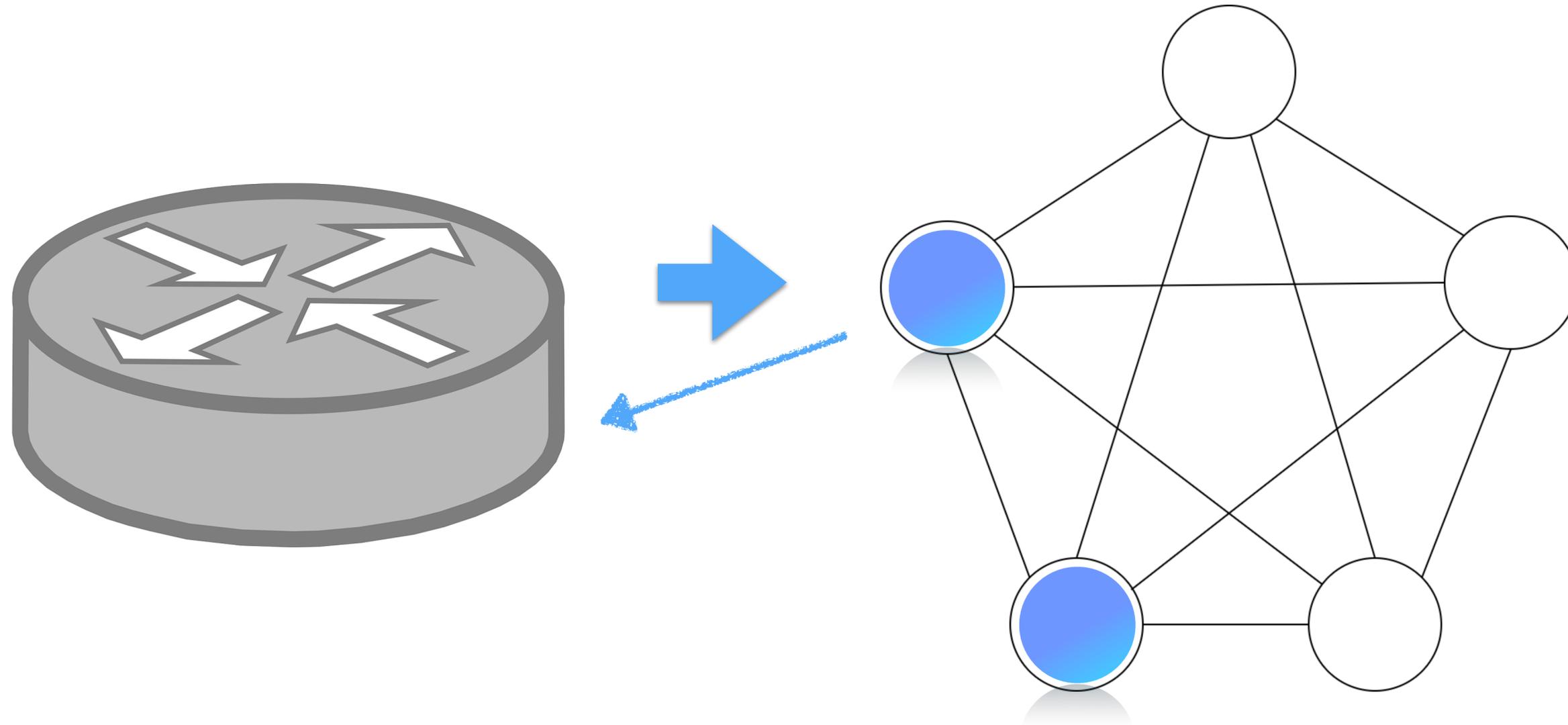
4.6 decentralized clusters - (kaa)

클러스터는 적어도 두벌 이상의 데이터가 존재



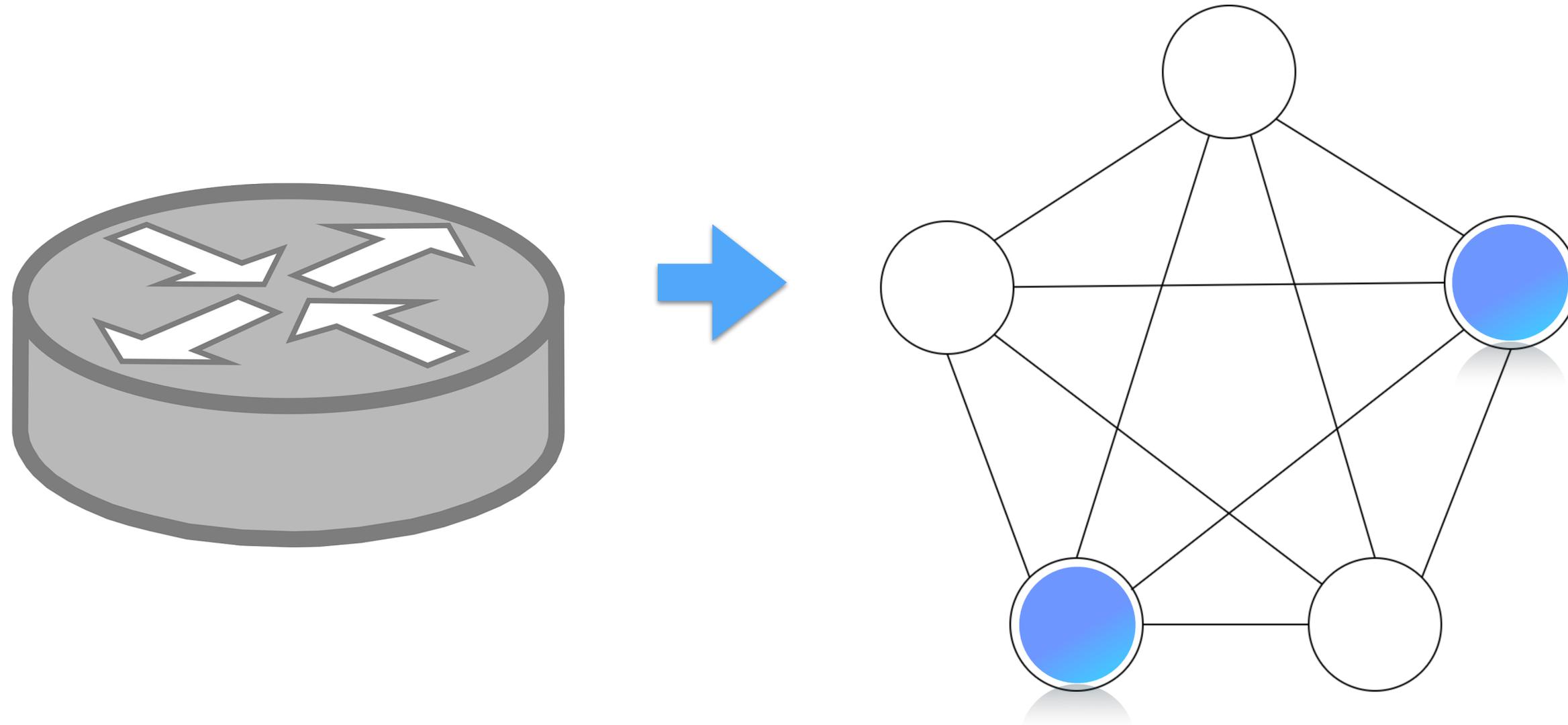
4.6 decentralized clusters - (kaa)

클러스터는 적어도 두벌 이상의 데이터가 존재



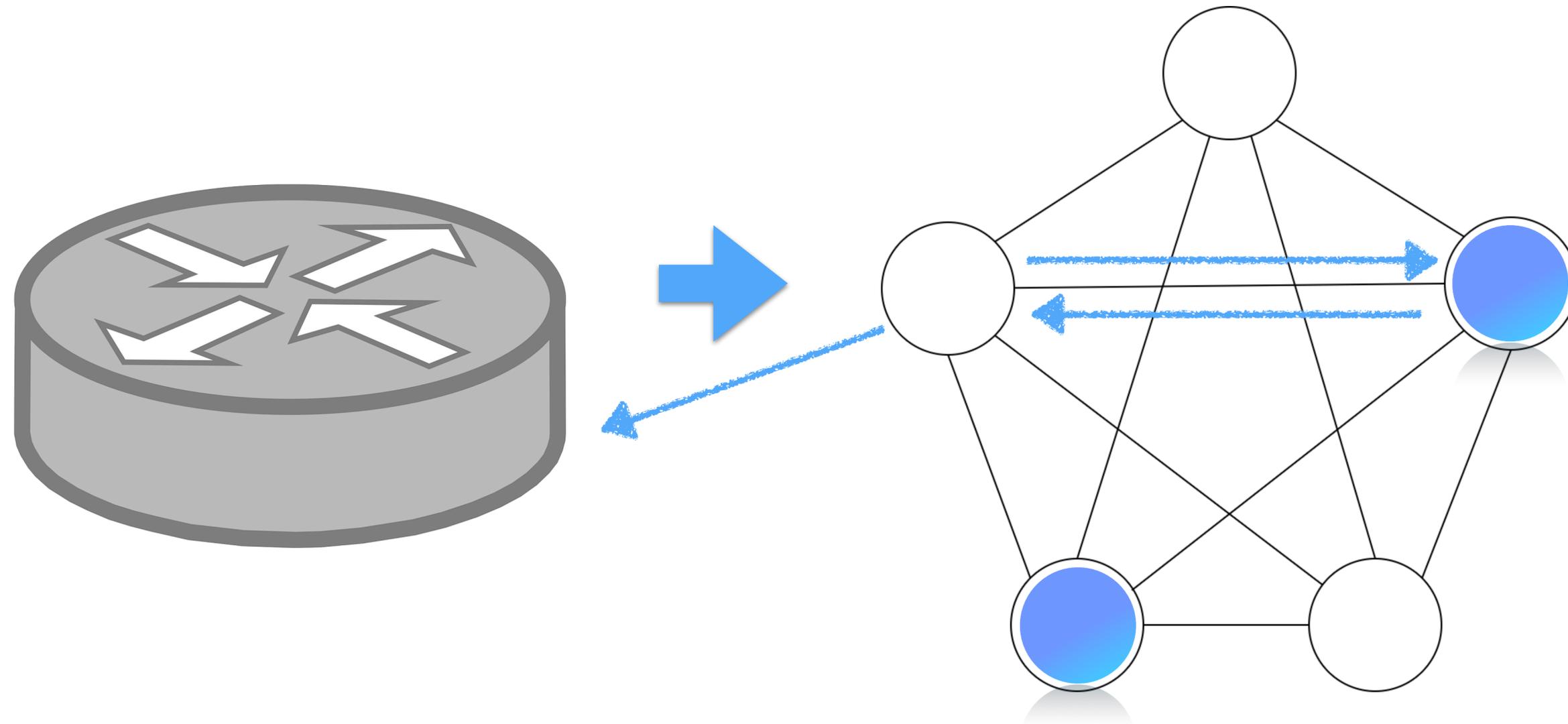
4.6 decentralized clusters - (kaa)

클러스터는 적어도 두벌 이상의 데이터가 존재



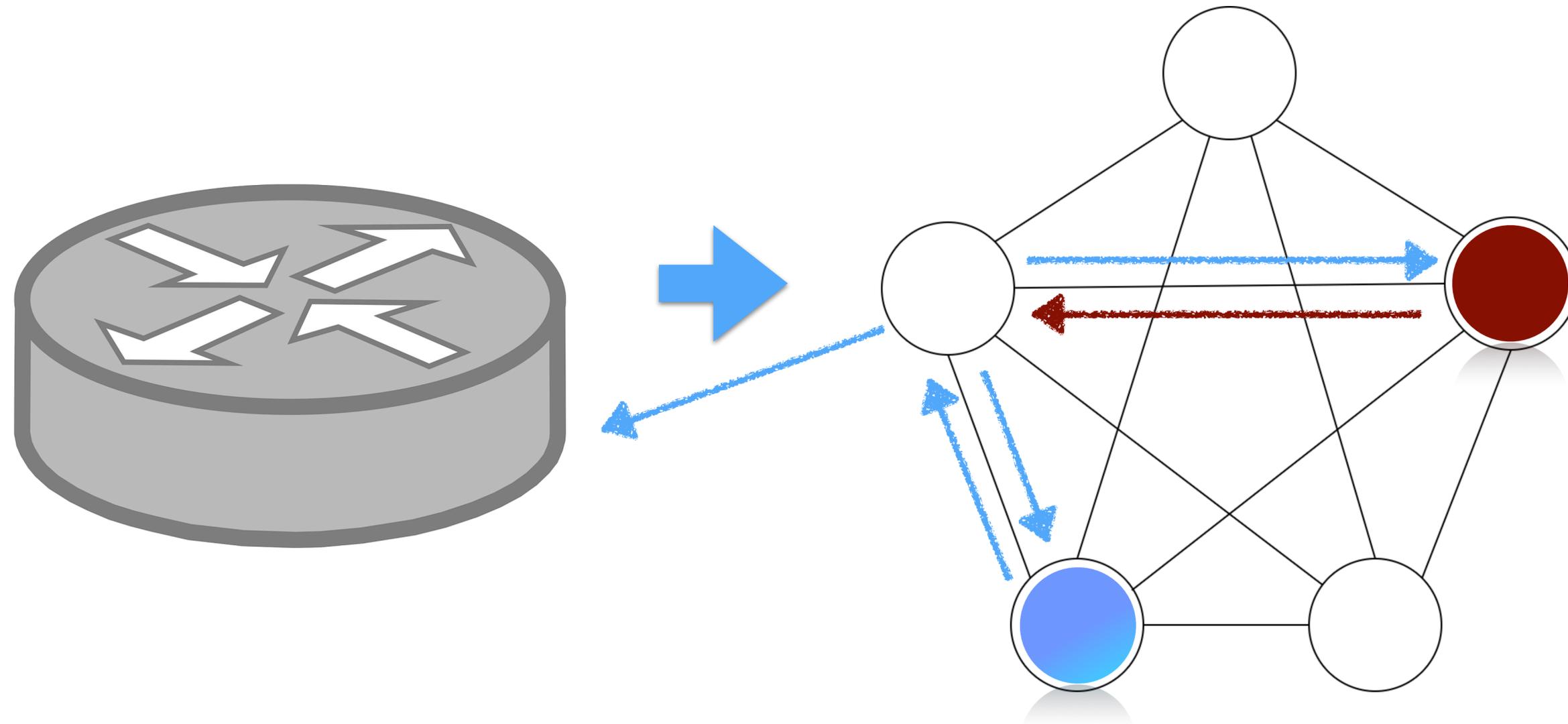
4.6 decentralized clusters - (kaa)

클러스터는 적어도 두벌 이상의 데이터가 존재



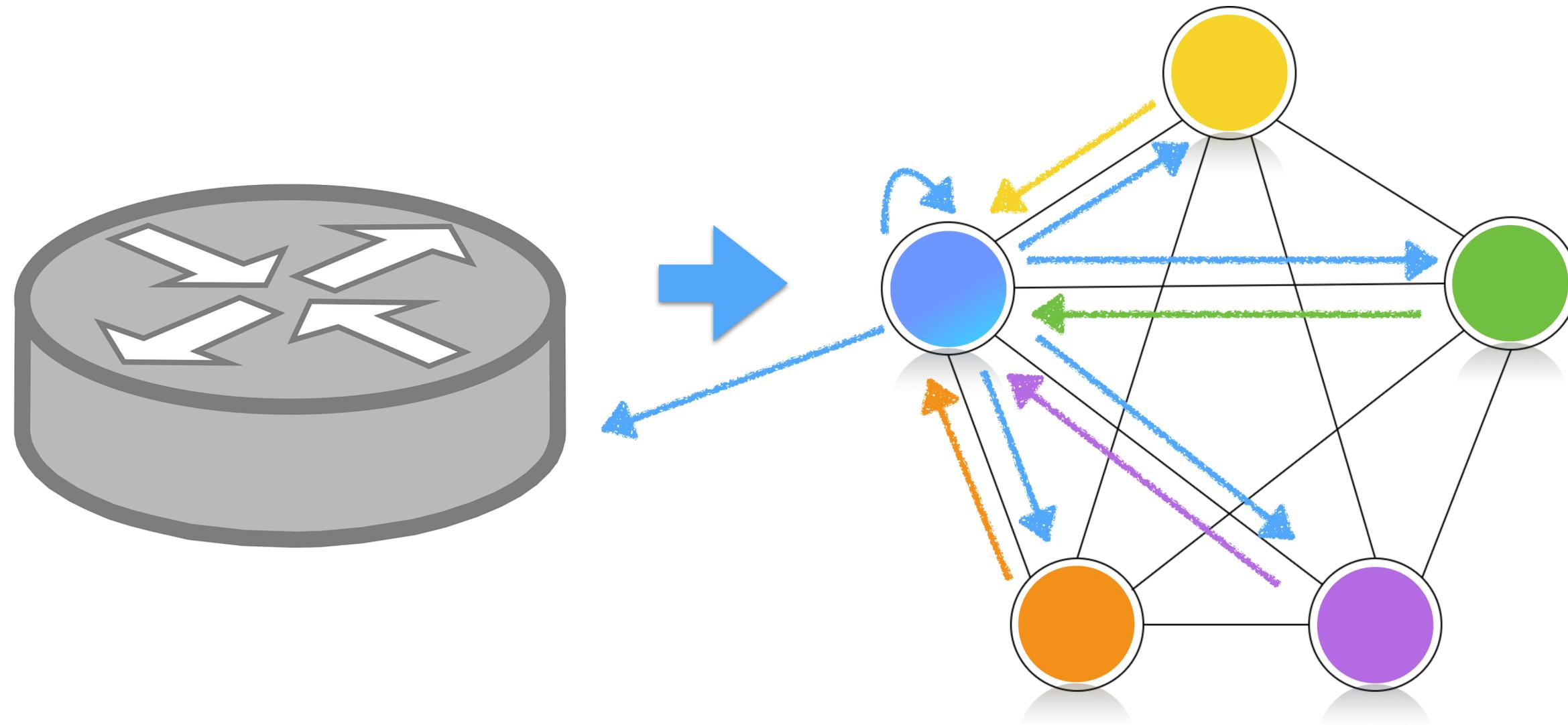
4.6 decentralized clusters - (kaa)

대상 노드가 죽거나 타임 아웃 경우 Dead Letter 처리



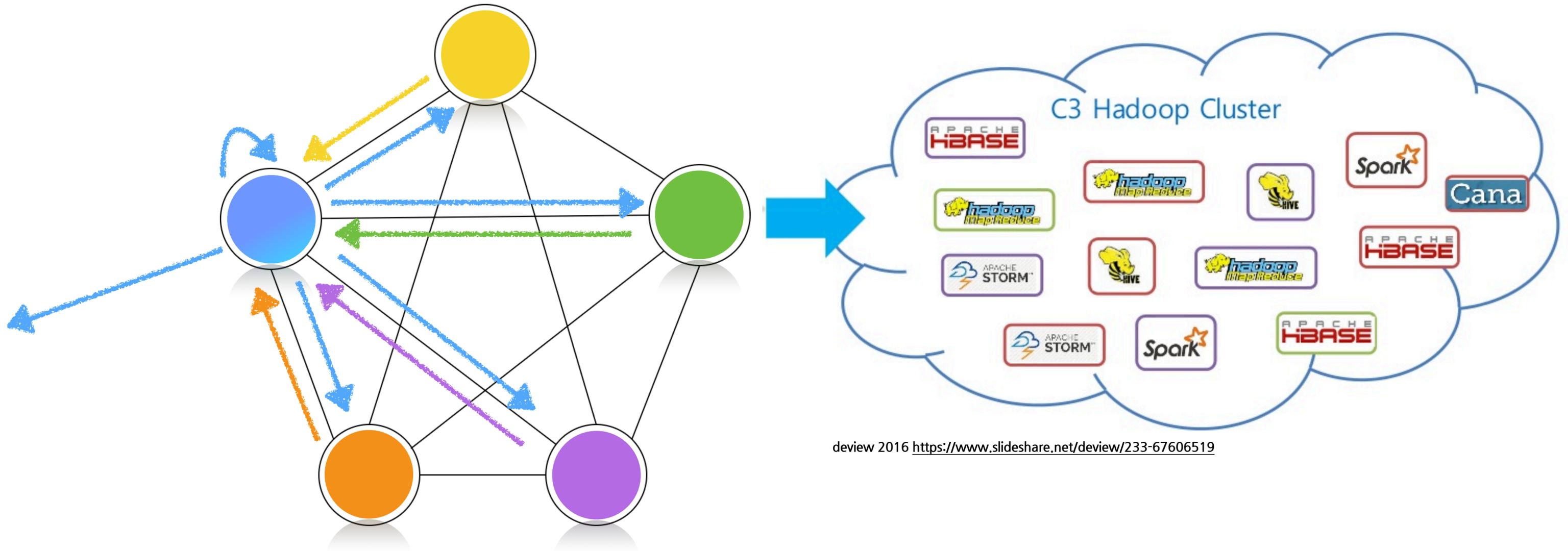
4.6 decentralized clusters - (kaa)

클러스터의 구성과 배포 옵션에 따라 성능, 안정성, 비용 고려



4.6 decentralized clusters - (kaa)

네이버의 검색을 지탱하는 대규모 인프라 시스템을 통한 모델 오프로딩



4.7 decentralized clusters - (kaa)

소수의 엔지니어로 15만개 이상의 모델을 학습하고 글로벌 환경에서 1만 여개의 모델을 서빙

CLOVA CHATBOT

사용자의 질문 의도를 이해하여 고객 대응 등 다양한 서비스에 활용할 수 있는 챗봇을 손쉽게 만들 수 있습니다.

- 네이버의 노하우가 축적된 자연어 처리 기술로 사용자의 질문 의도를 정확하게 파악하고 자연스러운 대화를 이어갑니다.
- 딥러닝 기반 자가 학습 알고리즘을 적용하여 발화 예시 입력 횟수를 줄이고 대화 모델을 최적화할 수 있습니다.

특징

독특한 챗봇 엔진

- 정교한 자연어 처리 기술과 머신러닝 기반 학습 알고리즘이 적용되어 자연어를 이해하고 빠른 답변 처리 가능

다국어 지원

- 한국어 뿐만 아니라 영어와 중국어, 일본어 지원
- 각 언어별로 최적화된 학습 알고리즘 적용

쉬운 챗봇 빌더

- 사용이 쉬운 챗봇 빌더를 통해 손쉽게 학습 데이터를 관리하고 발화 테스트 진행 가능

챗봇 FAQ

Q. 개발자 없이도 챗봇 제작이 가능한가요?

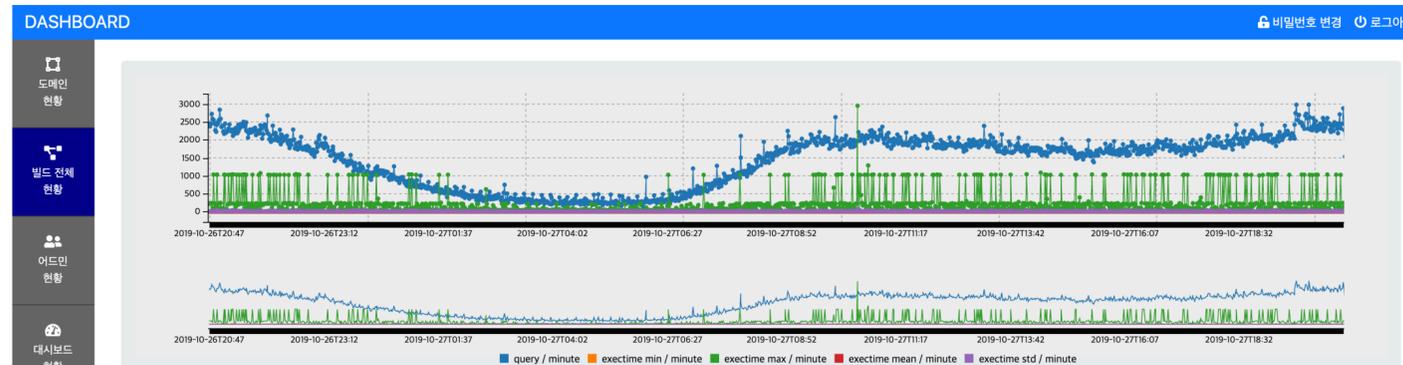
A. 네, 가능합니다. 제공되는 빌더를 통해 대화 모델 빌드 후 수 분에서 몇 시간 안에 학습 및 배포가 가능하며, 시뮬레이터를 통해 손쉽게 챗봇의 동작을 확인할 수 있습니다. 재학습을 위해 신규 발화를 추가로 등록하면 빌드 완료 후 새롭게 학습된 모델이 즉시 서비스됩니다.

Q. 음성으로 대답하는 챗봇을 제공하고 싶은데 어떻게 해야 하나요?

A. 네이버 클라우드 플랫폼의 다양한 API 상품을 활용해 챗봇 서비스를 확장할 수 있습니다. 음성 인식 및 합성 기술을 제공하는 Clova API와도 쉽게 연동해보세요. (Clova Speech Recognition / Clova Speech Synthesis)



Clova Chatbot
퀵스타트 가이드



cmgbldr-beta01	cceengine001	cceengine002	cceengine003	cceengine004	cceengine005
	4306 69548 TRAINING_C3 20:37:41.0	4281 69533 TRAINING_C3 20:18:57.0	4306 69551 TRAINING_C3 20:37:53.0	4306 69552 TRAINING_C3 20:37:59.0	4306 69548 TRAINING_COMPLETED 20:37:35.0
	4306 69558 TRAINING_C3 20:38:43.0	4306 69550 TRAINING_C3 20:37:48.0	4306 69561 TRAINING_C3 20:38:54.0	4306 69557 TRAINING_C3 20:38:32.0	4306 69558 TRAINING_C3 20:38:38.0
	4306 69564 TRAINING_C3 20:39:16.0	4306 69565 PREPARE_FOR_KAA 20:38:17.0	4306 69566 TRAINING_C3 20:39:28.0	4306 69562 TRAINING_C3 20:39:00.0	4306 69563 TRAINING_C3 20:39:06.0
	4306 69569 TRAINING_C3 20:39:44.0	4306 69560 TRAINING_C3 20:38:48.0	4306 69571 TRAINING_C3 20:40:04.0	4306 69567 TRAINING_C3 20:39:34.0	4306 69566 TRAINING_C3 20:39:39.0
	4306 69574 TRAINING_C3 20:40:24.0	4306 69565 TRAINING_C3 20:39:23.0	4306 69576 TRAINING_C3 20:40:37.0	4306 69572 TRAINING_C3 20:40:10.0	4306 69573 TRAINING_C3 20:40:18.0
	4306 69579 TRAINING_C3 20:40:54.0	4306 69576 TRAINING_C3 20:39:59.0	4306 69581 TRAINING_C3 20:41:06.0	4306 69577 TRAINING_C3 20:40:42.0	4306 69576 TRAINING_C3 20:40:49.0
	4306 69584 TRAINING_C3 20:41:32.0	4306 69570 TRAINING_C3 20:40:31.0	4306 69586 TRAINING_C3 20:42:09.0	4306 69582 TRAINING_C3 20:41:18.0	4306 69583 TRAINING_C3 20:41:27.0
	4306 69589 TRAINING_C3 20:42:30.0	4306 69575 TRAINING_C3 20:40:31.0	4306 69591 TRAINING_C3 20:42:41.0	4306 69587 TRAINING_C3 20:42:18.0	4306 69583 TRAINING_C3 20:41:27.0
	4306 69594 TRAINING_C3 20:42:58.0	4306 69580 TRAINING_C3 20:41:00.0	4306 69596 TRAINING_C3 20:43:13.0	4306 69582 TRAINING_C3 20:42:46.0	4306 69588 TRAINING_C3 20:42:25.0
	4306 69599 TRAINING_C3 20:43:30.0	4306 69585 TRAINING_C3 20:41:37.0	4306 69601 TRAINING_C3 20:43:41.0	4306 69597 TRAINING_C3 20:43:20.0	4306 69593 TRAINING_C3 20:42:52.0
	4306 69604 TRAINING_C3 20:44:03.0	4306 69580 TRAINING_C3 20:42:36.0	4306 69606 TRAINING_C3 20:44:14.0	4306 69602 TRAINING_C3 20:43:47.0	4306 69598 TRAINING_C3 20:43:25.0
	4306 69609 TRAINING_C3 20:44:33.0	4306 69595 TRAINING_C3 20:43:06.0	4306 69611 TRAINING_C3 20:44:46.0	4306 69607 TRAINING_C3 20:44:22.0	4306 69603 TRAINING_C3 20:43:53.0
	4306 69614 TRAINING_C3 20:45:02.0	4306 69600 TRAINING_C3 20:43:35.0	4306 69616 TRAINING_C3 20:45:19.0	4306 69612 TRAINING_C3 20:44:52.0	4306 69608 TRAINING_C3 20:44:28.0
	4306 69619 TRAINING_C3 20:45:36.0	4306 69605 TRAINING_C3 20:44:08.0	4306 69621 TRAINING_C3 20:45:47.0	4306 69617 TRAINING_C3 20:45:26.0	4306 69613 TRAINING_C3 20:44:57.0
		4306 69610 TRAINING_C3 20:44:40.0			4306 69618 TRAINING_C3 20:45:31.0
		4306 69615 TRAINING_C3 20:45:09.0			
		4306 69620 TRAINING_C3 20:45:42.0			

빌드 현황

- 최근 2000개의 빌드를 표시합니다. 클릭시 c3dl 빌드 현황 페이지로 이동

Copy CSV Excel PDF Print Show 10 entries Search:

빌드 ID	도메인 ID	생성일	C3DL ID	상태	DT	auto_ml	spark_appid	ETC
68103	2870	2019-10-26 15:15:31.0	1570432872011_40913	COMPLETED	2019-10-26T06-15-39.006Z			cceengine003-chatbot.nfra.io



Q & A

Thank You