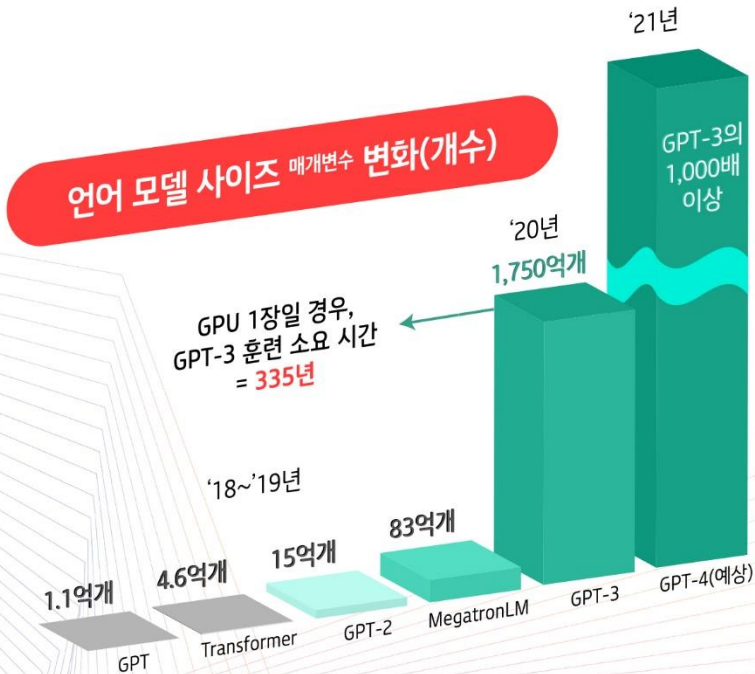


# Hyperscale AI Computing

## 서비스 소개

kt Cloud

## AI 산업 발전에 따른 GPU Computing의 급격한 수요 증가



- 전 산업에 걸친 AI Application 활용도 증가
- AI 모델 추론을 위한 대규모 고성능 인프라 필요

### N사, 검색 서비스 전체에 GPT-3와 같은 대규모 AI 모델 적용

※ 박성은 기자 © 입력 2021.05.08 12:10 ■ 댓글 0 ♡ 좋아요 0

AI-검색 연구자 대상 'N사 검색 콜로키움' 행사서 7일 발표  
슈퍼컴퓨팅 제품으로 '엔비디아 DGX 슈퍼팻' 사용 중  
대규모 생성모델 활용 관건인 개인정보보호, 연합학습으로 해결 예정

### 이통사·포털, '초거대 AI' 기술 개발 경쟁 치열

슈퍼컴퓨팅 활용, 무궁무진 확장성 특징  
원천 기술 확보 기반 '포스트 AI 시대' 선점 경쟁  
美 일론 머스크, 기존 대비 매개변수 17배 늘린 GPT-3 공개  
네이버·카카오·SKT·KT·LGU+, 대학 공동연구 및 세계 학회 논문 발표도

신희강 기자 입력 2021-12-02 05:56 | 수정 2021-12-02 09:02

특정 벤더, on-premise 위주로 공급되고 있으나 효율적인 운영이 어려운 상황

## On-premise 구매 방식

A100, DGX server, DGX superPOD 등



on-premise 구축 방식의 hurdles

- 1 GPU 서버 투자 비용 고가
- 2 딜리버리, 구축 등 공급에 시간 소요(6개월 이상)
- 3 자체 전산실 고전력 수용의 어려움
- 4 유희자원 활용률이 낮은 문제

※ N사 제품 위주

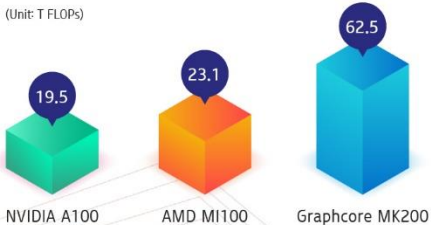
## N사 외 GPU 제품군

- 멀티 벤더 이기종 GPU 도입을 통해 원가 대비 높은 성능 구현 가능성 有
- 단 AI Framework의 Native 지원 부재



GPU 성능 비교(FP32 Peak 성능 기준)

(Unit: T FLOPs)



## Public Cloud 인스턴스 이용\*

CSP가 제공하는 GPU Server  
(T4, V100, A100...)

V100 GPU Memory 128GB 기준 월 요금

A*사	V100	p3.16xlarge	2,100만원
N*사	V100	V100 GPU server	1,234만원

※ Pass-through GPU VM

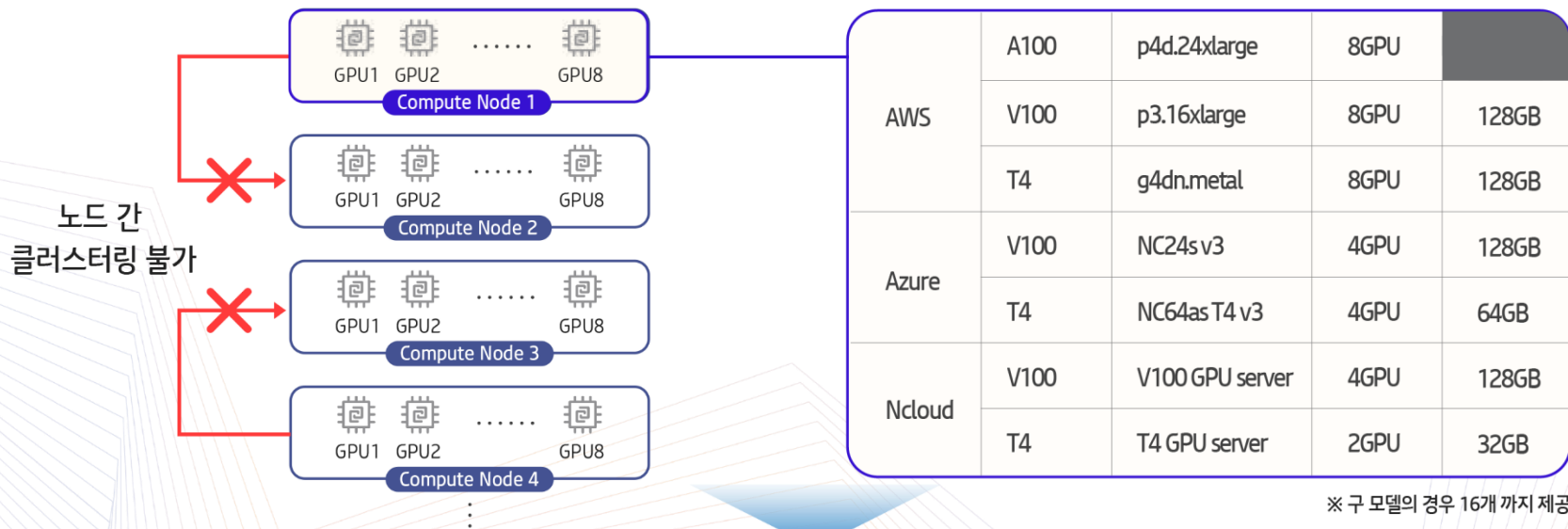
가성비 / 효율성을 갖춘 GPU 서비스 필요



### 1

### 대규모 GPU 클러스터링 지원 불가 (물리 노드 용량에 종속)

- 최대 8개 GPU 자원 제공
- GPU 카드를 VM에 귀속하는 Pass-through(Dedicated) 형태



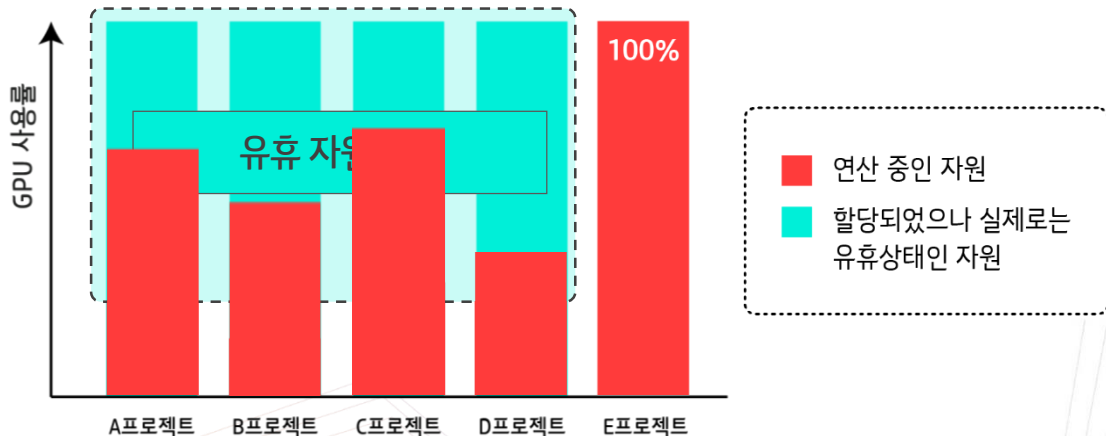
AI모델 대형화에 따라 대규모 GPU 인프라 필요

2

GPU 자원의 비효율적인 활용

연산 작업에 투입되지 않은 GPU를 활용할 수 없는 비효율성

GPU 자원의  
사용률 예시 ▶



공급기업의 구축비용 부담 가중, 수요기업의 이용료 부담으로 이어지는 문제

## 3

### 코드 수정의 불편함

#### 수정 고려 범위

- 각 GPU별 연산 분할
- 입력 데이터 분배 및 결과 취합
- 멀티 프로세스 실행 및 통신 초기화
- GPU 간 명시적 데이터 통신 추가
- 성능을 위한 병렬화 파라미터 조정

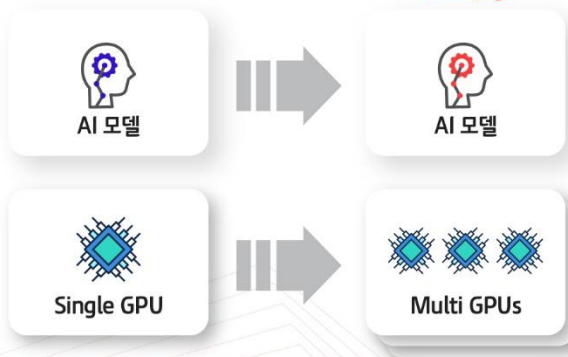
- 모델 고도화에 따라 Multi-GPU의 수요 증가
- Single GPU로 개발된 모델을 Multi-GPU로 변경할 경우, 사용자가 각 GPU에서 처리해야 할 연산을 할당하는 개발 변경 필요

```
class ToyModel(nn.Module):
    def __init__(self):
        super(ToyModel, self).__init__()
        self.net1 = nn.Linear(10, 10)
        self.relu = nn.ReLU()
        self.net2 = nn.Linear(10, 5)

    def forward(self, x):
        return self.net2(self.relu(self.net1(x)))

model = ToyModel().to('cuda')
output = model(input)
loss = loss_fn(output, target)
loss.backward()
optimizer.step()
```

#### 코드 수정



모델이 커질수록 개발자의 수동 배분 작업 및 연산 처리의 비효율 증가

```
class ToyModel(nn.Module):
    def __init__(self, dev0, dev1):
        super(ToyModel, self).__init__()
        self.dev0 = dev0
        self.dev1 = dev1
        self.net1 = nn.Linear(10, 10).to(dev0)
        self.relu = nn.ReLU()
        self.net2 = nn.Linear(10, 5).to(dev1)

    def forward(self, x):
        x = x.to(self.dev0)
        x = self.relu(self.net1(x))
        x = x.to(self.dev1)
        return self.net2(x)

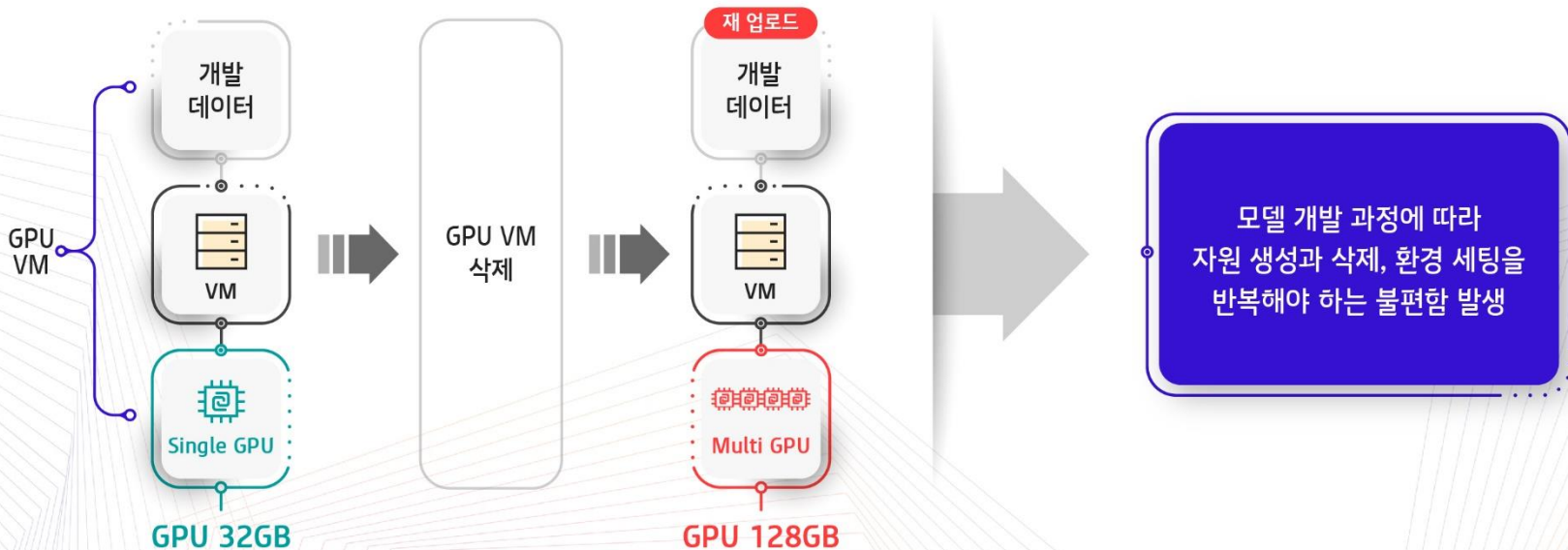
def execute(rank, world_size):
    os.environ['MASTER_ADDR'] = 'localhost'
    os.environ['MASTER_PORT'] = '12355'
    dist.init_process_group("gloo", rank=rank,
        world_size=world_size)
    dev0 = (rank * 2) % world_size
    dev1 = (rank * 2 + 1) % world_size
    model = ToyModel(dev0, dev1)
    ddp_model = DDP(model)
    output = ddp_model(input)
    loss = loss_fn(output, target)
    loss.backward()
    optimizer.step()
    dist.destroy_process_group()

n_gpus = torch.cuda.device_count()
mp.spawn(func, args=(n_gpus,), nprocs=n_gpus, join=True)
```

### 4

### 자원 규모 변경 시 확장성 부족

- 모델 개발 단계에 따라 필요한 연산의 규모는 상이
- GPU 자원 규모를 변경하기 위해 GPU VM 삭제 후 처음부터 재 생성 필요





## 03 Hyperscale AI Computing 서비스

- 20년부터 Moreh사와 협력하여 새로운 형태의 AI GPU 서비스 준비

### AI GPU 통합자원 관리 플랫폼

KT Cloud Portal / Console

KT Cloud DX 플랫폼 (Openstack)

GPU 관리

AI Framework  
PYTORCH  
TensorFlow

SD GPU Library  
MoDA  
moDNN

가상 자원 관리

Compute Network Storage Identity

KVM Hypervisor



200G 초고속 네트워크 (Infiniband 스위치)

AMD



AMD



Lustre  
File System

AI용 병렬처리  
파일 시스템

12. 10 상용 출시

ktCloud G-Cloud → 라이브! 상품 솔루션 고객지원 파트너 프로모션 로그인 🔍 클라우드 콘솔 ⓘ

**GRAND OPEN**

# 국내최초! 대규모 클러스터링이 가능한 완전한 종량제 GPU 서비스

## Hyperscale AI Computing

### 출시기념 무료 EVENT

KT의 가상화 기술을 통해 GPU의 대규모 클러스터링을 지원하고, 실 사용량에 기반하여 합리적인 종량제 요금으로 GPU 자원을 제공하는 최초의 클라우드 서비스입니다. 22년 2월 28일까지 무료, GPU 이용료에 한정하여 무상

상품 보기 →    프로모션 보기 →



## 1

### 대규모 GPU 자원 활용 가능

모델 대형화에 대응할 수 있는 수백~수천 개의 대규모 GPU 클러스터링 지원

기존

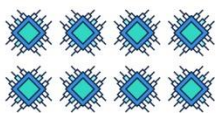
클러스터링 규모 제약으로 인한 AI 모델 대형화 한계

AI서비스

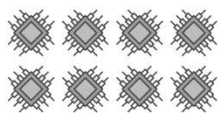


개별 물리 장비에 할당된 GPU수량 이상으로 확장 불가

Compute node 1



Compute node 2



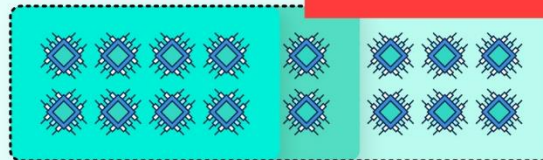
KT

AI 모델 대형화

AI서비스



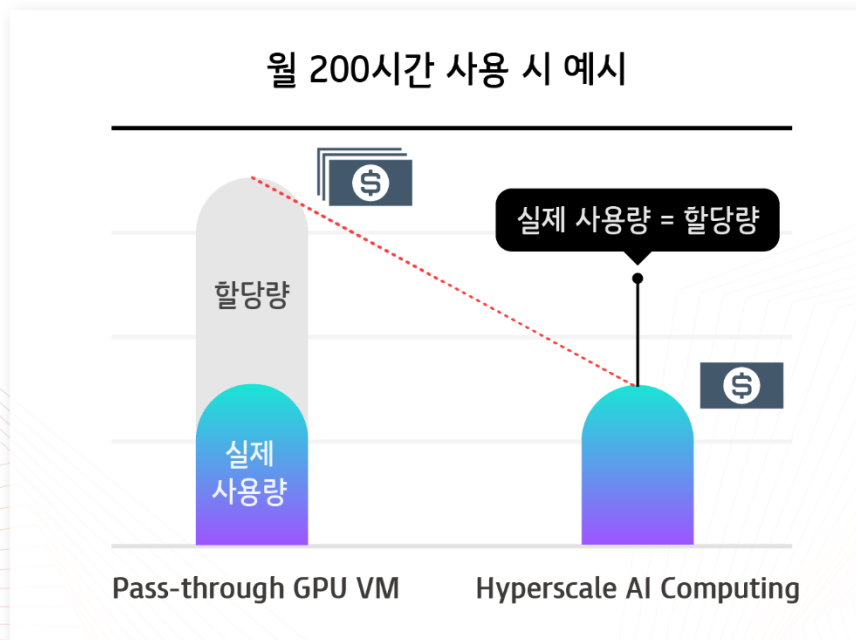
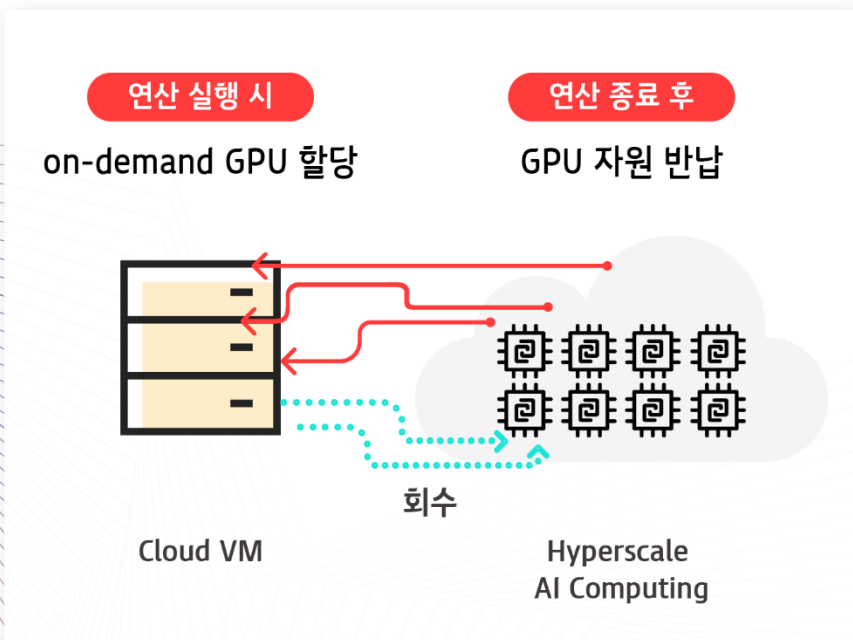
원하는 만큼 GPU Cluster 확장



## 2

### 실 사용량 기반의 종량제 서비스

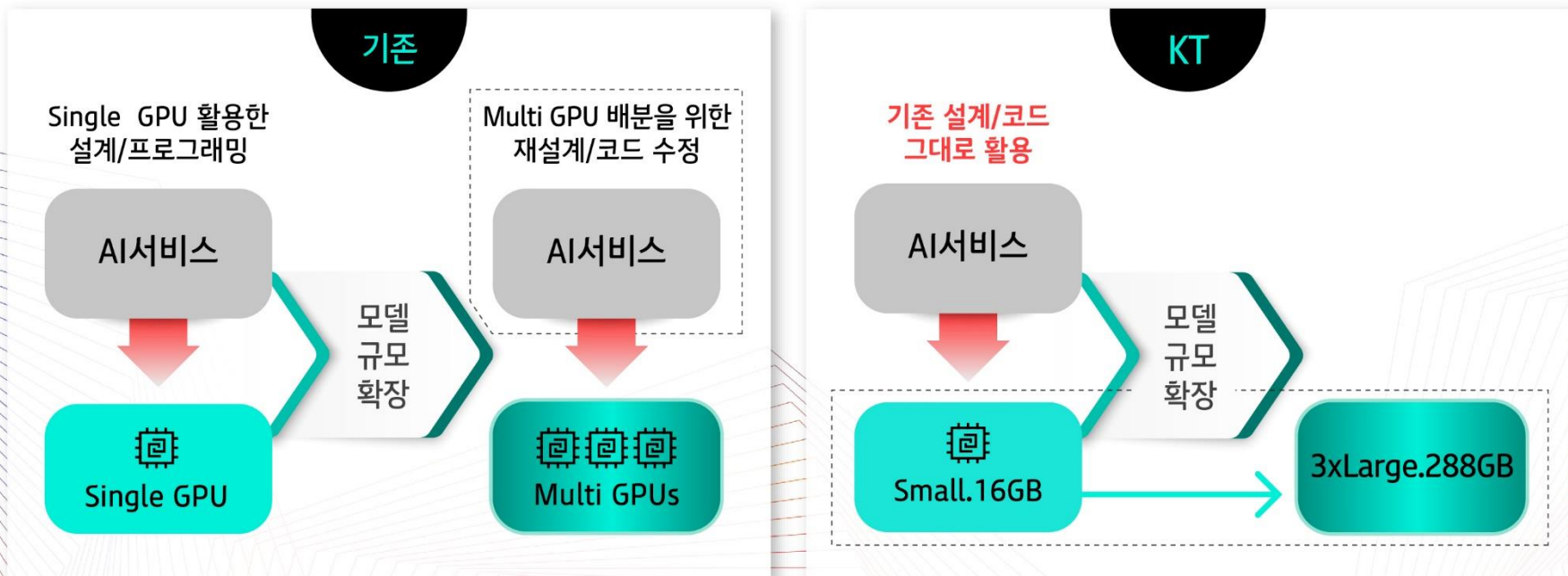
최초의 사용량 기반의 종량제 GPU 서비스  
연산 실행 시에만 GPU 자원을 할당·과금하여, 사용 패턴에 맞추어 합리적으로 비용 절감 가능



## 3

### 모델 프로그래밍 호환성

별도 개발 없이 Hyperscale AI Computing 내 컴파일러가 자동으로 GPU 자원 분산 처리 가능

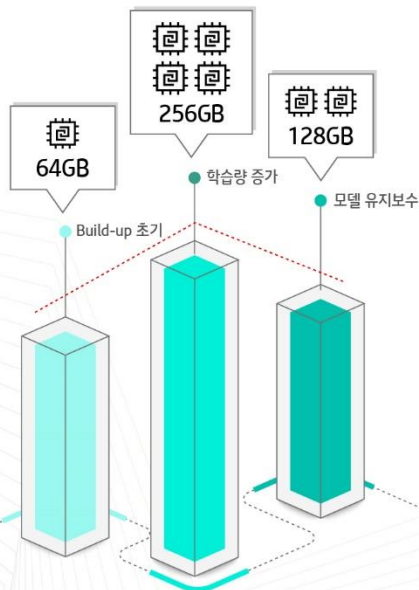




## 4

### 개발의 유연성 및 연속성 제고

AI 모델 개발 스테이지에 따라 필요한 자원 규모를 자원 Live 상태에서 탄력적으로 변경 가능



▶ 개발 단계별 필요 GPU 메모리 규모 예시

### Server

가상 서버를 관리합니다.

서버 생성   시작   정지   재시작   삭제   접속설정   ...   모든 위치 · 모든 상

<input type="checkbox"/>	이름 ↓	상태	위치
<input checked="" type="checkbox"/>	GPU-0407	● 사용	DX
<input type="checkbox"/>	hm-dnd-watch-test	● 사용	DX
<input type="checkbox"/>	kjh-test33	● 사용	DX
<input type="checkbox"/>	Inrssi0902	● 사용	DX
<input type="checkbox"/>	InrssiServer2	● 사용	DX
<input type="checkbox"/>	Inrtest-VM	● 정지	
<input type="checkbox"/>	Inrtest0804	● 정지	
<input type="checkbox"/>	nhsSSD	● 정지	

- 상세정보
- 사양변경
- AI가속기 사양변경**
- 비밀번호 확인
- OS 초기화
- 추가 사실IP
- Volume 관리
- 이미지 생성
- 요금제 변경
- D1 플랫폼 해지

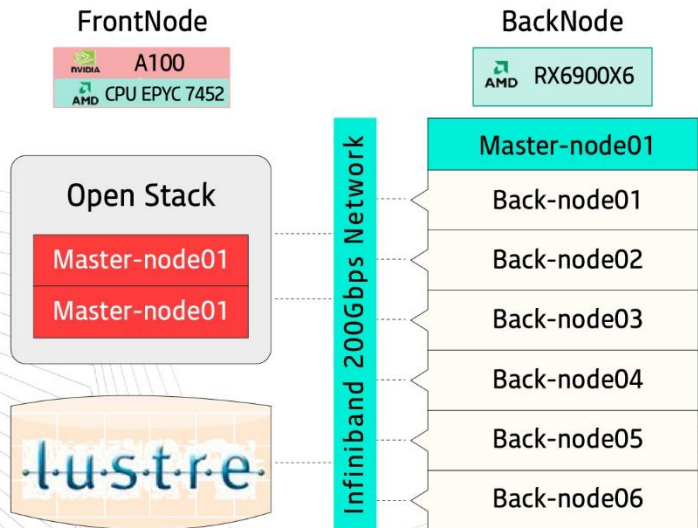
**Small.16GB**

- medium.32GB
- large.48GB
- large.64GB
- xlarge.96GB
- 2xlarge.192GB
- 3xlarge.288GB

# 03 Hyperscale AI Computing 서비스

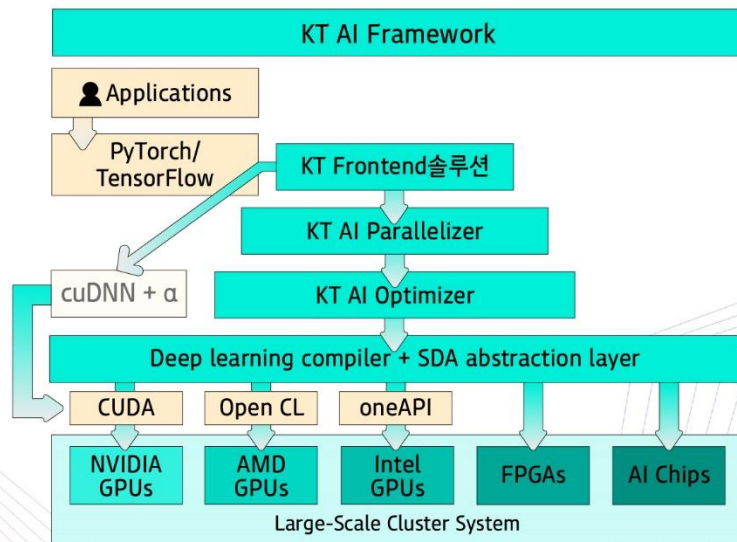
- Public Cloud 제공을 위한 H/W, S/W 최적화 설계 제작 및 개발 추진

## H/W 구성도



- ▶ Customer용 GPU를 우선 활용하여 12월 서비스 개시, 향후 데이터센터용 GPU 적용 예정(~'22)
- ▶ 공조를 위한 별도 냉방 장치 고안 등 맞춤형 장비 설계

## S/W 구성도



- ▶ 기존 코드의 100% 유지가 가능한 사용성 확보
- ▶ NVIDIA 외 이종 GPU를 동일하게 활용할 수 있는 Hyperscale AI Computing의 AI Framework 개발

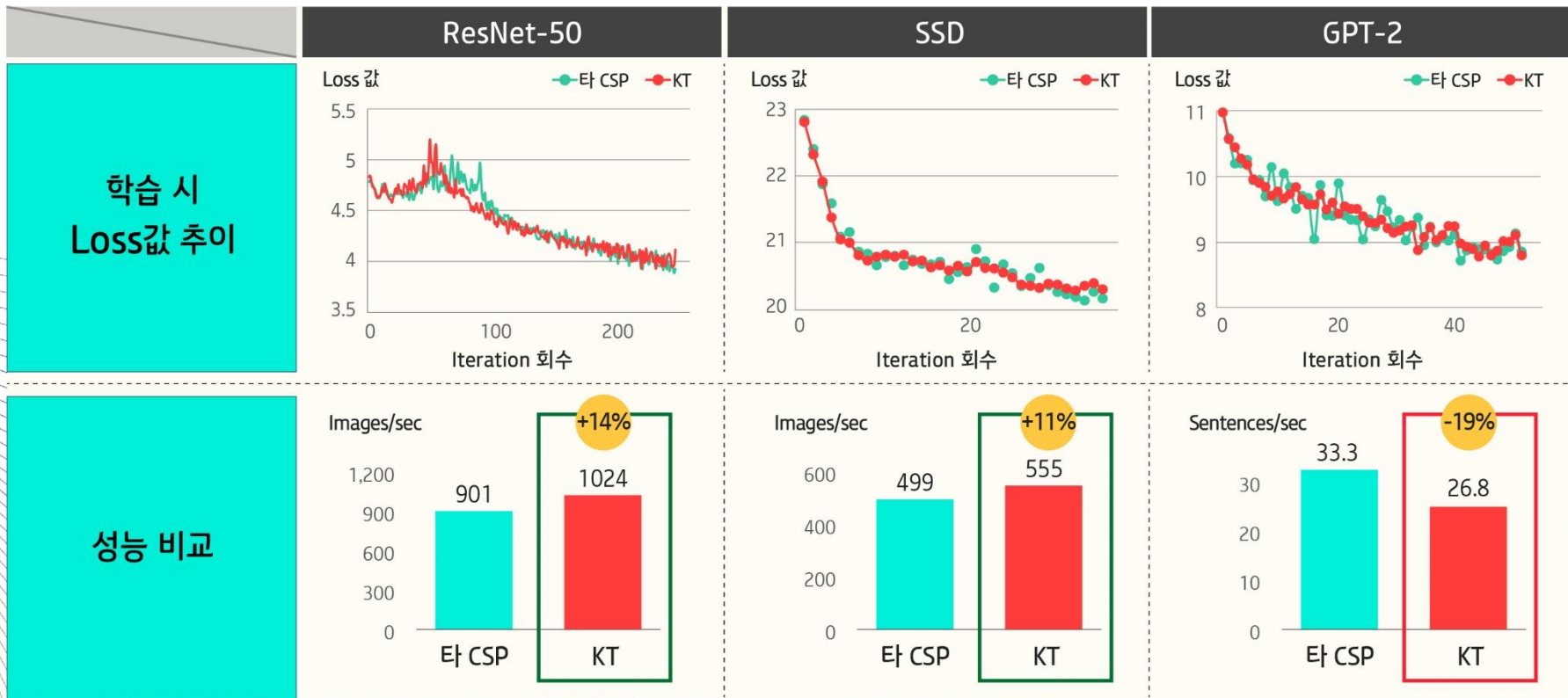
## 04 Hyperscale AI Computing 성능 검증

항 목	실험 대상	테스트 내용
Single Node 정합성/성능 평가	KT Hyperscale AI Computing 타 CSP의 V100 GPU Server 비교	<b>Benchmark 모델에 대해</b> (ResNet-50, SSD, GPT-2) <b>정합성</b> 학습 진행에 따른 loss값 감소 (모델 정확도 증가) 추이 비교 <b>성능</b> 초당 학습량 비교 <small>※ ResNet-50: 대표적인 사물 분류 모델 ※ SSD: 사물의 종류와 위치를 식별하는 모델 ※ GPT-2: 대표적인 Transformer 언어 모델</small>
Multi Node 성능 평가	KT Hyperscale AI Computing 에서 GPU 개수에 따른 성능 측정 4 → 8 → 16 → 24 → 32 1 node    2 node    3 node    4 node	GPU 수 증가에 따른 초당 학습량 향상 추이 측정 (GPT-2)



# 04 Hyperscale AI Computing 성능 검증

## 4-1. Single Node 정합성/성능 평가 결과

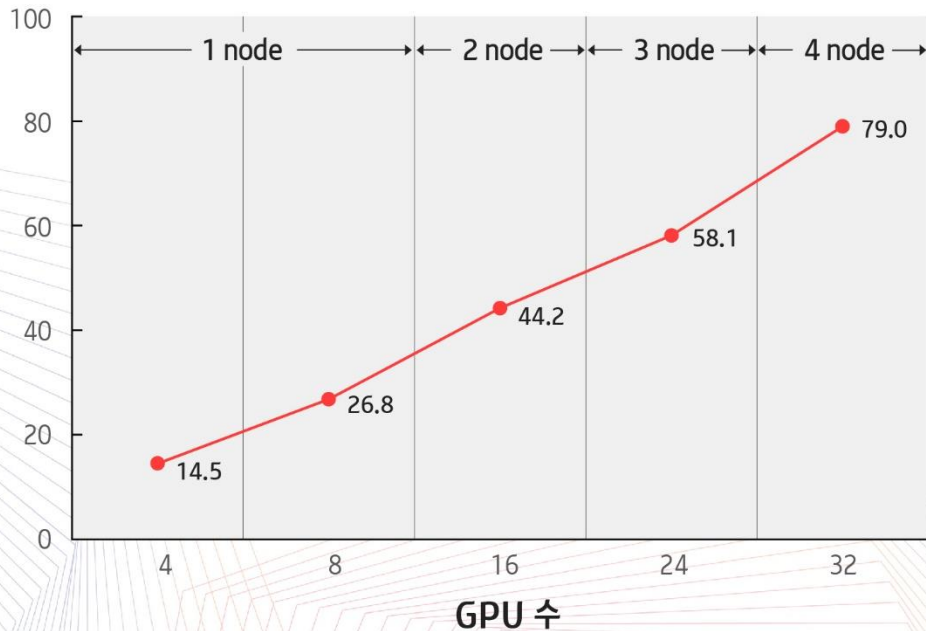


## 04 Hyperscale AI Computing 성능 검증

### 4-2. Multi Node 성능 평가

GPU/Node 구성 별 GPT-2 성능 지표

Sentences/sec



물리 node 간 클러스터링  
제약을 극복하고 GPU를  
사용할 수 있는 유일한 서비스

GPU/노드 개수가  
증가함에 따라 계산 성능이  
scalable하게 증가

1

16개 모델 제공, 매월 추가 → 41개 모델 순차 제공 (~1Q)

Vision	ResNet	SSD
	VGG	SqueezeNet
	Inception V3	Mask R-CNN
	ArcFace	3D U-Net
NLP	GPT-2	Transformer
	BERT	RNN-T
	Tacotron	
기타	DLRM	
	NCF	

+

월 단위  
모델 추가

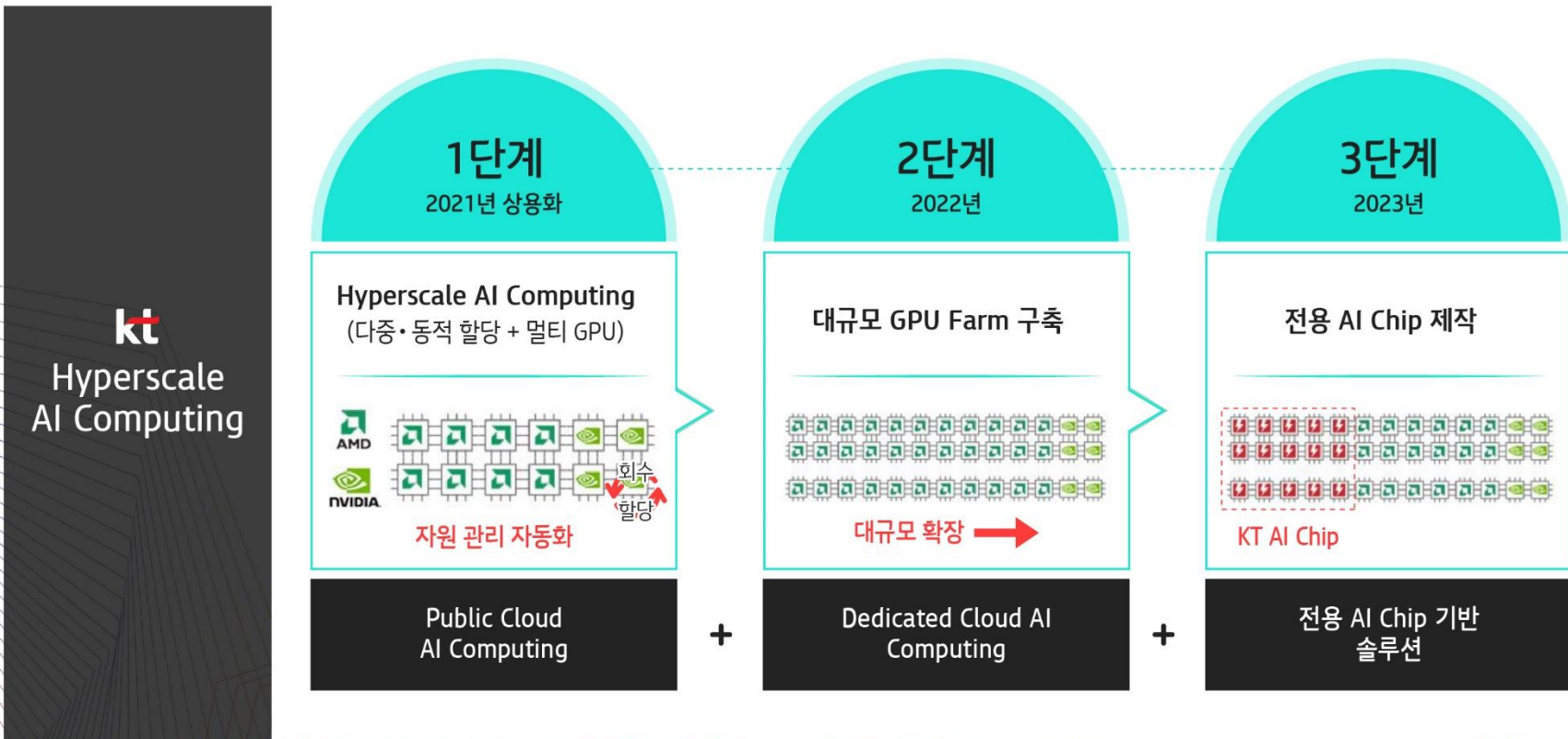
2

Tensorflow 확장 제공 (~2Q)

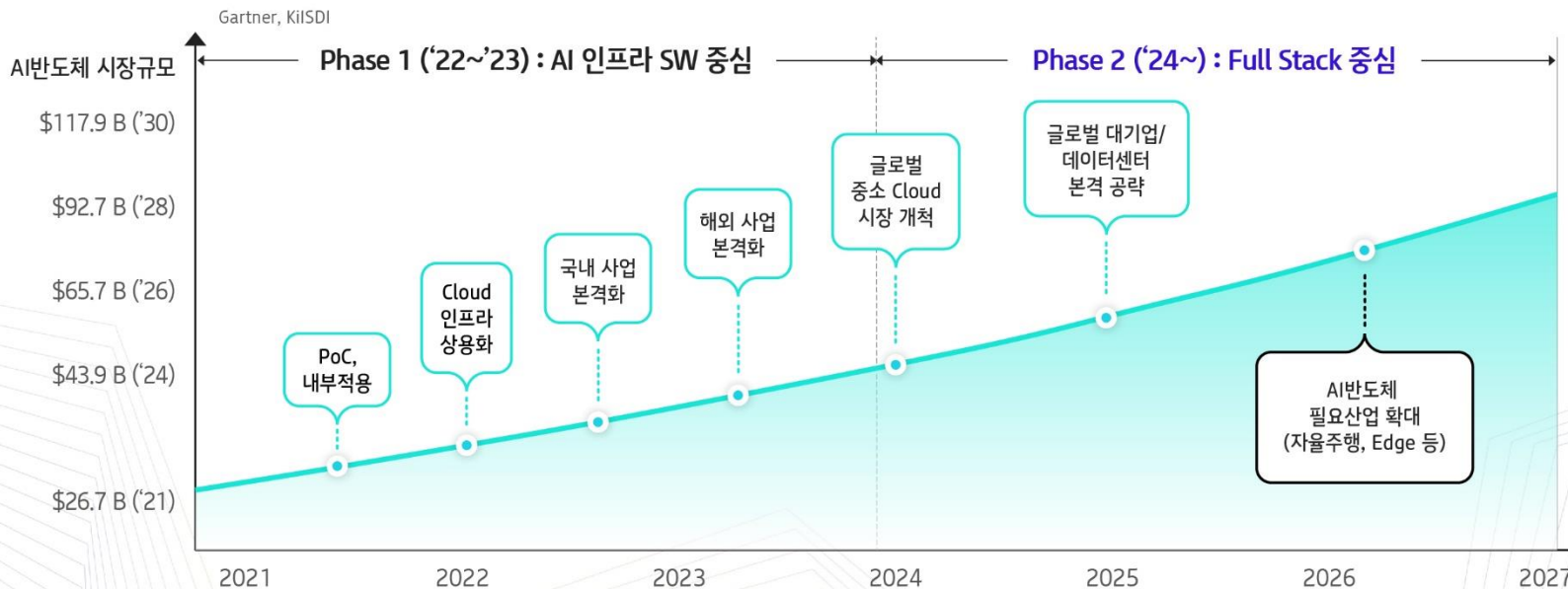




순수 국산 기술로 AI Full Stack 영역까지 서비스 확대(~23년)



## 향후 로드맵 : AI Full Stack 중심 글로벌 확장



주요 Initiative	파트너십 구축	SW부터 시작하여 시장 개척	Full Stack으로 전면 확장	글로벌 AI 리더십 구축
		<ul style="list-style-type: none"> <li>KT-Moreh SW PoC 및 상용화</li> <li>양사 전략적 제휴 추진</li> </ul>	<ul style="list-style-type: none"> <li>외부 NPU(AMD)기반 SW 단독사업화</li> <li>AI반도체 자체 개발(F사에서 spin off)</li> </ul>	<ul style="list-style-type: none"> <li>AI반도체+SW의 Full Stack 사업화</li> <li>선진국 시장에서 본격 경쟁</li> </ul>

KT는 고객의 “AI서비스 실현”을 위한  
튼튼한 조력자(Enabler)입니다.



kt Cloud