



Company Introduction

Make Al accessible for all.

Company Name	TEN Inc.
CEO	Sejin Oh
Established	July 3, 2020
Address	No. 1203, Hyeonik Building, Taehaeran-ro, Yeoksam-dong, Gangnam-gu, Seoul, South Korea
Main Services	MLOps solution, 'Al Pub' Al infrastructure consulting service 'RA:X'
Employees	22 employees
Website & Social Media	<u>https://ten1010.io/</u> <u>https://www.youtube.com/@ten1300</u> <u>https://www.linkedin.com/company/ten1010/</u>

Since launching our services in 2020, we have been selected for various projects and awards, recognized for our strong competitiveness and potentials.

2021

- 02 Started business by launching the AI Pub service07 Established TEN Inc.
- 03 Joined NVIDIA Inception Program
- 04 Received seed funding from Kingsley Ventures
- 05 Official partnership with NetApp
- 06 Acquired Venture Business Certification Selected for N&UP Program¹
- 07 Selected for TIPS² Startup Growth and Technology Development Project
- 08 Established company research center
 Received K-Startup Award in the "Al Service" category at the Korea
 Leading Company Awards
- Attended NVIDIA GTC³ selected as Top Startup from Korea
 Selected as service provider for AI Voucher Project

- 01 Received Award of Excellence on 2021 Global Scale Up IR Day of Global Business Partnership Program
 - 03 Selected for Al Voucher Project
 - 07 Received K-Startup Award in the "MLOps Platform" category at the Korea Leading Company Awards
 - 08 Selected for Initial Startup Package of Korea Institute of Startup & Entrepreneurship Development
 - 11 Certified as Company with Technological Excellence by KoData (T4 Level)
 - 12 Received pre-series A round funding (Ascendo Ventures, Quantum Ventures Korea, Korea Credit Guarantee Fund)

0000
シロノミ
2023

2024

- 01 Registered 5 patents related to resource management (KOR) (10-2488614, 10-2488615, 10-2488618, 10-2488619, 10-2488620)
- 02 Hosted the First Al Workshop⁴
- 07 Signed NetApp Preferred Partnership
- 08 Signed IBM Partnership
- 09 Signed Redhat Partnership
- 11 Completed project on high-density GPU-based cluster computing system for deep learning computation for Korea Agency for Defense Development
- AI Pub : GS Certification Grade 1 from TTA'MOST Promising Korean Tech Company 2023' awarded by CIO Review
- 01 Industry-Academic Cooperation Agreement on AI Infrastructure with Sungkyunkwan University SAL < COMPASS LAB
- 04 AI Alliance MOU with Kolon Benit
- 05 Awarded the 19th Digital Innovation Grand Prize for IT

- and Ministry of SMEs and Startups aimed at supporting businesses to enter the global market
- TIPS: Tech Incubator Program for Startup
- 3) NVIDIA GTC: NVIDIA GPU Technology Conference



¹⁾ N&UP Program: Joint project between NVIDIA



committed to making Al more accessible^{Mission} for all by creating tools that lower the barrier to technology.^{Vision}



In a world where Al is more accessible, everyone can create value and enjoy the benefits of using Al.

TEN offers affordable AI tools that are easy to use to make AI more accessible to all. We envision a world where experience and knowledge related to AI is shared throughout the entire society.



How are you preparing for Al in the era of deep learning?

We are living in the era of deep learning, where AI is at the center of attention of all members of society, from individuals to businesses. In the very near future, everyone will be able to convert their knowledge into data and automate it through training.

Unlike the times of statistical pattern recognition and physical phenomenon modeling, today, the ideas of the individual creating the AI have become ever more important.

Statistical Pattern Recognition

HMMs for Speech

• Example of using HMM for word "yes" on an utterance:



Deep Learning



Building a Cat Detector using Convolutional Neural Networks — TensorFlow for Hackers (Part III) | by Venelin Valkov | Medium



"We are at the iPhone moment of Al." - Jensen Huang

Startups are competing to create innovative products and business models, while existing companies explore ways to respond.



GTC 2023, Jensen Huang (NVIDIA CEO)



We are currently in the process of reaching AI singularity.

More than ever before, these times require talents, data, and infrastructure to make AI more widely available.



TEN

We are gradually achieving the conditions required to enable the general public, rather than just certain groups of experts, to build their own AI and train them using their knowledge converted into data.

However, the issue of **cost** still remains when it comes to **infrastructure**, and it cannot be solved with time and effort alone.



Since the emergence of deep learning, the performance of AI has been predominantly determined by computing power. Unlike traditional statistical pattern recognition, computing power has a significant impact on the performance for deep learning.

To enhance performance, significant investment is required. That is why resource efficiency will become increasingly important going forward.

Statistical Pattern Recognition

Deep Learning

Performance ∞ Computing power

Weak proportional relationship



Strong proportional relationship



Companies of all sizes and industries are striving to build high-performance server infrastructure.

The linear scale clearly shows the exponential increase in the required computing resource. This is a clear indication of the imminent Al singularity and the importance of infrastructure.



TEN

However, improving the efficiency of infrastructure resources is a challenging task, as the continuous increase in resources required for AI models makes it more complicated to build an AI infrastructure.

With more GPU machines required to cover the increased computing demand of AI models, building an AI infrastructure has become more complex and challenging.





High efficiency in infra utilization : the key to making AI more widely available.







TEN offers MLOps software and dedicated Al infrastructure that facilitate the development and operation of Al.



TEN provides solutions for each infrastructure management point to address issues that arise during AI model training and operation.

Problems in Al Training Process

- Operation platform optimized to the model training characteristics (resource-hungry) is needed when building on-premise infrastructure.
- As resource is allocated in server units with the conventional platform, it is not possible to know the idle and operation rates at the GPU unit level.
- Overhead is incurred in the resource allocation process to change the server settings (OS, drivers, libraries, etc.) for each developer.
- Job scheduling is required for each model training.

Solution for AI model training

Maximized GPU resource operation rate 24/7

Problems in AI Operations Process

- GPU resource specification and amount required for service operation cannot be estimated for model deployment.
- Stability cannot be secured when operating multiple services due to the interference between services within the GPU resource.
- Cost goes up from the continuous increase in the amount of resources used by model.

Solution for AI model operation

Stable operation quality with minimum GPU usage

TEN offers COASTER, a container platform, and Al Pub Dev and Al Pub Ops, which are MLOps tools and also RA:X, an Al Infrastructure consulting service.

Container platform

GUI-based MLOps service for AI development and operations

For AI Infrastructure

RA:X

TEN's Products

Al training and operation require numerous skilled experts across various fields. TEN developed function-level products that support companies to establish their own governance structures and offers MLOps tools to operate Al services without all-rounder experts.

TEN

TEN's **COASTER** and **AI Pub** provide optimum solutions that encompass everything from AI service to the building and maintenance of infrastructure.

1 Day

Shortened time to reach Al distribution, which used to take 8.6 months

100% util & 1/10 cost

Infrastructure is used at the highest operation rate during the AI training phase, then brought to the lowest rate during the AI operation phase.

Off-the-shelf

Ready-to-use platform without any complicated development process that requires skilled experts from various fields

Tailored

Build AI-tailored infrastructure that maximizes the performance of costly GPU resources

COASTER and **AI Pub** use Kubernetes, which is the most popular platform for container orchestration.

10 INSIGHTS ON REAL-WORLD CONTAINER USE

FACT 2

Usage of GPU-based compute on Containerized workloads has increased

We observed a 58 percent year-over-year increase in the compute time used by containerized GPU-based instances (compared to a 25 percent increase in noncontainerized GPU-based compute time over the same time period).

TEN listened to the needs of our clients during AI development and developed the optimum solution with AI Pub.

 We purchased servers in bulk, but we are struggling with their operation due to the significant management overhead required to allocate them to different teams. *II* Teams purchased their own servers. Now we have different types of GPU servers that are difficult to manage together and the usage rate is very low.

 The engineering department purchased a GPU server for research work, but we do not have an effective solution to share it between multiple labs.

We purchased A100 GPU recently, but we are finding it difficult to configure the MIG function every time and share it with our developers.

II Multiple users are sharing the server, but we are unsure of how to operate spaces like storage or image registry to manage the data of each user.

We manage our development products using docker images. It is difficult to manage the images that are used by each user and the images that have to be shared with the team by the administrator.

TEN listened to the needs of our clients during AI operation and developed the optimum solution with AI Pub.

II We developed a speech synthesizer system, but we lack the experience to launch the service on a GPU server.

11

We developed an Al service, but we are struggling with stable operation due to a lack of operators and knowhow to operate the GPU server.

We do not have a solution to manage the various types of services and different docker images for each service version.

II The server resource has to be divided and managed by each service operator or team, but we do not know how. We developed an AI service for identifying product defects. Is there software that can operate and manage it automatically?

||

11

II We need to monitor the traffic of each service and check and manage response failures. Is there an effective way to do this?

11

Al Pub is a fully-managed solution that focuses on addressing issues related to the operation and infrastructure utilization throughout the Al lifecycle.

It is designed to provide service to any type of AI model in various infrastructure environments.

		AI Models	s/Services		
Vision	Speech	Chatbot	Text	Search	Recommend
		P AlPub -	MLOps Platform		
Train	Deploy	Manage	Monitor	Orchestration	UI
	C	COASTER - GPL	J integration and fr	agmentation tech	nnology
		GPU Server I	nfrastructure		
On-premise	e Pr	rivate Cloud	Public Clou	d H	ybrid Cloud

We offer all-in-one Al-specific services in partnership with vendors specialized in Al.

Enterprises and universities in various areas are using TEN's services.

SUNG KYUN KWAN UNIVERSITY

26

Our aim was to create an AI deployment tool that can be widely used. We envisioned a friendly and welcoming "pub-like" environment, where individuals can come together to develop and operate AI.

Al Public Al Publich Al Publish

COASTER

Container platform with functions improved by extending Kubernetes

COASTER is a container platform that extends Kubernetes and enhances the operation of GPU infrastructure and user management functions.

Main Services	Service Description
GPU fragmentation	Fragment the utilization and memory of one GPU unit into 100 blocks
Inquire and allocate GPU resource	Inquire the computing resource in the entire cluster using extended Kubernetes commands
User entitlement management – by individuals and groups	Assign policies to users and groups
Job scheduling and priority management	Kubernetes-based job scheduler automatically initiates jobs, while al so allowing a human operator to manually re-prioritize them through GUI.

FUNCTION of COASTER 1. COASTER supports the usage of fragmented GPUs.

Allocating GPUs to containers in block units not only allows multiple containers to run on a single GPU but also ensures stability by preventing the interference of resource usage between containers.

Kubernetes Native : GPU Allocation

GPUs can only be allocated to containers as a whole unit, and multiple containers cannot run on a single GPU.

Coaster Extended : GPU Allocation

The utilization and memory of a single unit of GPU can be fragmented into 100 blocks, each spanning 1% increments.

FUNCTION of COASTER 2. COASTER enables the inquiring and allocation of GPU resources.

With the extended Kubernetes commands, you can inquire the computing resource of the entire cluster and allocate the required number of blocks for the appropriate type of GPU to each container. This user experience is different from the basic Kubernetes environment, where you can only inquire resource in server units and are unaware of the GPU type of each server.

Kubernetes Native : View GPU resource

: 39
4
83873772Ki
0
0
16093900Ki
1

Can only view the resource status of accessed node

Kubernetes Native : Allocate GPU resource

Allocate entire GPU

Coaster Extended : View GPU resource

NAME	RESOURCE_NAME	TOTAI	FRFF
block-tesla-t4	ten1010.īo/block-tesla-t4	400	130
сри	сри	16000m	9200m
gpu-tesla-t4	ten1010.io/gpu-tesla-t4	4	2
memory [root@aws-master	memory -01 ~]#	100001/238K1	1283834K

Can view the computing resource of the entire cluster using extended commands

Coaster Extended : Allocate GPU resource

Allocate divided GPU

FUNCTION of COASTER 3. COASTER provides user entitlement management on individual and group level.

Users who share a policy can be placed in the same group, enabling collective management of their entitlement to access resources, such as namespace, shared storage, image registry and server node.

Kubernetes Native : user entitlement management

The access to resources of each user is managed separately, making the management process complicated and difficult.

Coaster Extended : user entitlement management

Users who share a policy are in the same group, and their entitlement to access resources, such as namespace, shared storage, image registry and server node, can be managed together.

Coaster's open source project : Resource-group-controller

Example for user entitlement management of COASTER

Group A : User entitlement Only for User 1, 2, 3 Server node a / Name space a / Image registry a / Shared storage a

Group B : User entitlement Only for User 2, 5 Server node b / Name space b / Image registry b / Shared storage b

Group C : User entitlement Only for User 4, 6, 8, 9 Server node c / Name space c / Image registry c / Shared storage c

FUNCTION of COASTER 4. COASTER's scheduler makes it easy to re-prioritize jobs in the queue.

Kubernetes Native : scheduler

Basic scheduler of Kubernetes manages

Coaster Extended : scheduler

It is easy to change the priority of the jobs

See a demonstration of **COASTER**'s functions in this video.

System Info:						
Machine ID:	3d5c05376530a2eb49e3e90576f83c5b					
System UUID:	EC2C8B43-79BC-0A59-0CEB-008C83663BDF					
Boot ID:	31cfc8f9-155d-44bb-ac1b-f07f034b38b0					
Kernel Version:	3.10.0-1062.12.1.el7.x86_64					
OS Image:	CentOS Linux 7 (Core)					
Operating System:	linux					
Architecture:	amd64					
Container Runtime Version:	docker://23.0.1					
Kubelet Version:	v1.21.1					
Kube-Proxy Version:	v1.21.1					
PodCIDR:	10.244.0.0/24					
PodCIDRs	10.244.0.0/24					
Non-terminated Pods:	(8 in total)	CDU C	CDU Linite	N	Name and Links	
Namespace	Name	CPU RE Wests	CPU LIMITS	Memory Requests	Memory Limits	Age
luba flannal	kuha floppol da p7x2p	100-12	2E0m (C%)	100M÷ (0%)	EEOM: (2%)	11
kube system	Rube-Tlannel-us-p/x2p		250m (6%)	TOOMI (0%)	170M; (1%)	1.2m
kuba-system	coredns - 558bd4d5db - Nyp8n	10m (2	0 (0%)	70Mi (0%)	170Mi (1%)	12m
kube-system	etcd_aws_master_01	100m (2	0 (0%)	100M; (0%)	0 (0%)	12m
kube-system	kube-aniserver-aws-master-01	250m (**)	0 (0%)	0 (0%)	0 (0%)	12m
kube-system	kube-controller-manager-aws-mater-01	200 (5%)	0 (0%)	0 (0%)	0 (0%)	12m
kube-system	kube-proxy-zmb92	(0%)	0 (0%)	0 (0%)	0 (0%)	12m
kube-system	kube-scheduler-aws-master-01	100m (2%)	0 (0%)	0 (0%)	0 (0%)	12m
Allocated resources:						
(Total limits may be over 1	.00 percent, i.e., overcommitted.)					
Resource Requests	Limits					
cpu 950m (23	1%) 250m (6%)					
memory 340Mi (2	1%) 890Mi (5%)					
ephemeral-storage 100Mi (0	9%) 0 (0%)					
hugepages-1Gi 0 (0%)	0 (0%)					
hugepages-2Mi 0 (0%)	0 (0%)					
Events:						
Type Reason	Age From Message					

AIPub Dev

With COASTER at its core, fully-managed service that supports AI development

With COASTER at its core, AI Pub Dev supports AI development.

It provides a fully-managed service for model training, resource and workload management, etc.

Al Pub Dev

Main services	Service description
	Manage the user's development environment in docker image
Create workspace	Create workspace based on development images
	Connect with Jupyter Notebook and TensorBoard
Model training	Automatically allocate resource required for each AI training
	Can apply for GPU resource and CPU resource
	Limit resource usage of each user account
Descurse menogement	Withdraw idle resource
Resource management	Manage workspace of each node & Set MIG for each node
	Monitor entire infrastructure
Posourco group management	(for manager) Create resource group and Set user's authorization
	Edit resource group
Workload management	Stop/resume scheduler
	Job scheduling and prioritization
Lloogo biotony monogoment	Manage resource usage history of each user account
Usage history management	Download usage history

ten

FUNCTION of AI Pub Dev 1. AI Pub Dev can allocate AI infrastructure to a team or an individual user. It can also measure the amount of resource usage of each team or developer.

With AI Pub Dev, you can build an AI development infrastructure to allocate and manage resource during centralized management. You can also measure the resource usage of each user account.

FUNCTION of AI Pub Dev 2.

Al Pub Dev allows users to create workspaces based on resources accessible within the affiliated group.

TEN

See a demonstration of AI Pub Dev's functions in this video.

AIPubOps

With COASTER at its core, fully-managed service that supports AI operation

Al Pub Ops offers technologies required for Al operation.

Built on COASTER with Kubernetes, it provides service mesh and application functions.

Kubernetes Coaster

Inquire, allocate and manage all cluster resources based on fragmented GPU units

Allocate fragmented GPUs to containers and improve efficiency of required resource usage per service

Al Pub Ops supports Al operation with COASTER at its core. It provides a fully-managed service for Al service and resource.

Al Pub Ops

Main services	Service description
	Provides UI for service creation, suspension, deletion and distribution
Service creation and update	Provides UI for non-stop service update
	Version management and service rollback
Convice monitoring	Monitor operation status using service list and details
Service monitoring	Service error alert and troubleshooting by checking logs
Resource group management	Enables administrator to create resource group and set user entitlement
Descurrent	Enables allocation of GPU blocks for each service
Resource management	Monitor real-time operation rate of GPU blocks and servers

FUNCTION of AI Pub Ops 1.

Al Pub Ops manages the user's services using docker images It reduces cost by fragmenting the GPU into the smallest size of 100 units.

Al Pub Ops enables non-developers to create, suspend, delete, update and rollback services.

TEN

FUNCTION of AI Pub Ops 2.

Al Pub Ops also supports Al service operation on on-premise servers and offers the same service operation and management functions that are provided in the public cloud environment.

Al Pub Ops can reduce your Al operation cost by up to 90%.

Instead of operating your AI service on a public cloud, you can reduce cost by building a server of the same size and using AI Pub. Your service can be operated with only 10% of the server resource by using our GPU fragmentation function, and you can also save cost by eliminating the manpower required to develop, maintain or repair functions for service operation.

[Infrastructure Cost for Five-Year Operation of 50 AI Services of Client Company]

(Actual example: 10 servers with 4 units of T4 managed by 4 mid-level developers \rightarrow 1 server with 4 units of T4 managed by Al Pub and no managing personnel)

TEN also offers **consulting services** to address the difficulties in Al operation that may not be fully covered by Al Pub Ops.

We operate your service for you using our in-depth AI operation knowhow.

With our services, you can focus solely on creating more AI values without spending extra time and money to maintain continuous service stability.

Al Pub Ops integrates with NetApp Astra and supports the data management, backup, migration and rollback of applications that consist of multiple micro-services.

<u>'Easier management of AI services and data on MLOps platform in Kubernetes</u> <u>environment'</u>

Al Pub Ops uses NeuVector's scanning function to filter docker images from various sources to ensure security.

Al Pub Ops can check for vulnerabilities in container images, nodes that form the clusters, and running images or jobs in a container.

COASTER X 👧 SUSE

CVE Scanning

The CVE system manages unique identifiers assigned to various vulnerabilities and system defects that may expose devices, systems and programs to hacking.

Al Pub Ops provides monitoring services for the Al resource manager and service operator.

Use AI Pub Ops to monitor the status of the system and manage risks.

(³ AIPub									
통한 알렘 전티	서비	스 목록 8개							
서비스 성성	- 1-1								
AND THE									세비수 영정
A DECEMBER OF	No	리소스 그를 (서비스 분류 :	서비스 아름 0	운영기간 :	포트링보	리소스 타입	세블리카 개수 (소대 🕈
이미지 유해	-	기반 리소스 그룹	tts	voice-man002	2024/02/27 11:26:05 - (운영용)	844.5.3464.5138 URM: 80/984: 32222	testa-64-7% 17	03/178	898880 (8487)[83]
	2	기본 리소스 그룹	tta	voice-man001	2024/02/27 11:12:29 - (운영용)	294942:346645118 489:8080/999:31112	tesia-64-1%; 20 78	578/578	8427
	3	기본 리소스 그룹	im	chat-inbound	2024/02/27 11:37:31 ~ (記句書)	월4주소: 34645138 내부: 100/ 위부: 31113 (- 프로알림	tosia-64:17	178/178	주역 <mark>()</mark> (8세일기) 로그
	4	기번 리소스 그를	Im	chat-inbound-main	2024/02/27 12:17:27 - (世俗音)	월수주소: 34,64,5138 189(: - / 919): - 고도상장	tesia-64-1%; 6.78	578/578	[###2][#3]
	5	기본 리소스 그룹	lim	chat-promo	2024/02/27 23:20:36 - (2018)	월수주소: 34.64.5118 대학: 900/의학: 32010 고프 같은	tesia-64-1%; 5.78	571/571	(8487) (83)
	6	기번 리소스 그룹	voice	voice-man-003	2024/02/28 13:39:50 - (899)	84442 34645118 1899 811/1997 32213 1864/8	tesia-14-1%; 5.78	2.11/2/1	(SVEN) ED
		利用利止から意	.10	vsioi-weman001		12月1日-11-12日 12月1-1-12日 12月1-1-12日 12月1-1-12日 12月1-12日	1000-01,278	0.11/11	생성 대기용 성서보기 물고
		7日 月五六 2番	10	voloe-woman001		20474:3464518 (201-/109- (26.98)	teriis-14, 2 M	0.31/121	생성 대기용 [상세모기] 루그]
		7世 町本本 口道	m	voior-woman001		D++ha: 3464.5138 GP(-/ 598) - [15:49]	tusin-14, 2.78	0.78/129	상성 대기용 실서M기 문고

See a demonstration of AI Pub Ops' functions in this video.

AI Reference Architecture of TEN

TEN's know-how and test-based consulting service for AI infrastructure construction **RA:X** does not suggest that each hardware be purchased on a budget. **RA:X** makes you know the value of the entire infrastructure.

TEN introduces a **Reference Architecture** built and operated directly.

Infrastructure configuration consulting based on RA:X follows the process outlined below.

TEN

RA:X

TEN is hosting seminars and workshops on the challenges and solutions of Al infrastructure configuration. We empathize with concerns about infrastructure and aim to share our insights in various forums to assist in finding solutions.

AI 인프라와 NVIDIA 인증 솔루션 통합 구성 워크샵

나에게 딱 맞는 AI 인프라 조합을 찾는 여정

RA:X Nov. 23 Thu 14:00~17:00 드림플러스 강남

2023 AI Infrastructure Workshops : New AI Infra Solution of TEN & NVIDIA-Certified Solution Opening Session by TEN

2023 AI Infrastructure Workshops : New AI Infra Solution of TEN & NVIDIA-Certified Solution Session C, Closing Session by TEN

Are you looking for a competent partner for your AI services?

Experience a whole new level of AI development and operation by partnering with TEN Inc. We offer you highly-efficient solutions as your trustworthy partner.

- A No. 1203, Hyeonik Building, Taehaeran-ro, Yeoksam-dong, Gangnam-gu, Seoul, South Korea
- **P** +82-2-6956-1071
- M <u>helloten@ten1010.io</u>
- H https://ten1010.io/

Make Al accessible for all

Copyright © 2024 Ten Inc. All rights reserved