

Reinforcement Learning

나는 강화학습으로 축구한다

Google Research와 Manchester City F.C.의
인공지능 축구 프로젝트를 대한민국에서
재조명하고 직접 구현해보는 캠프



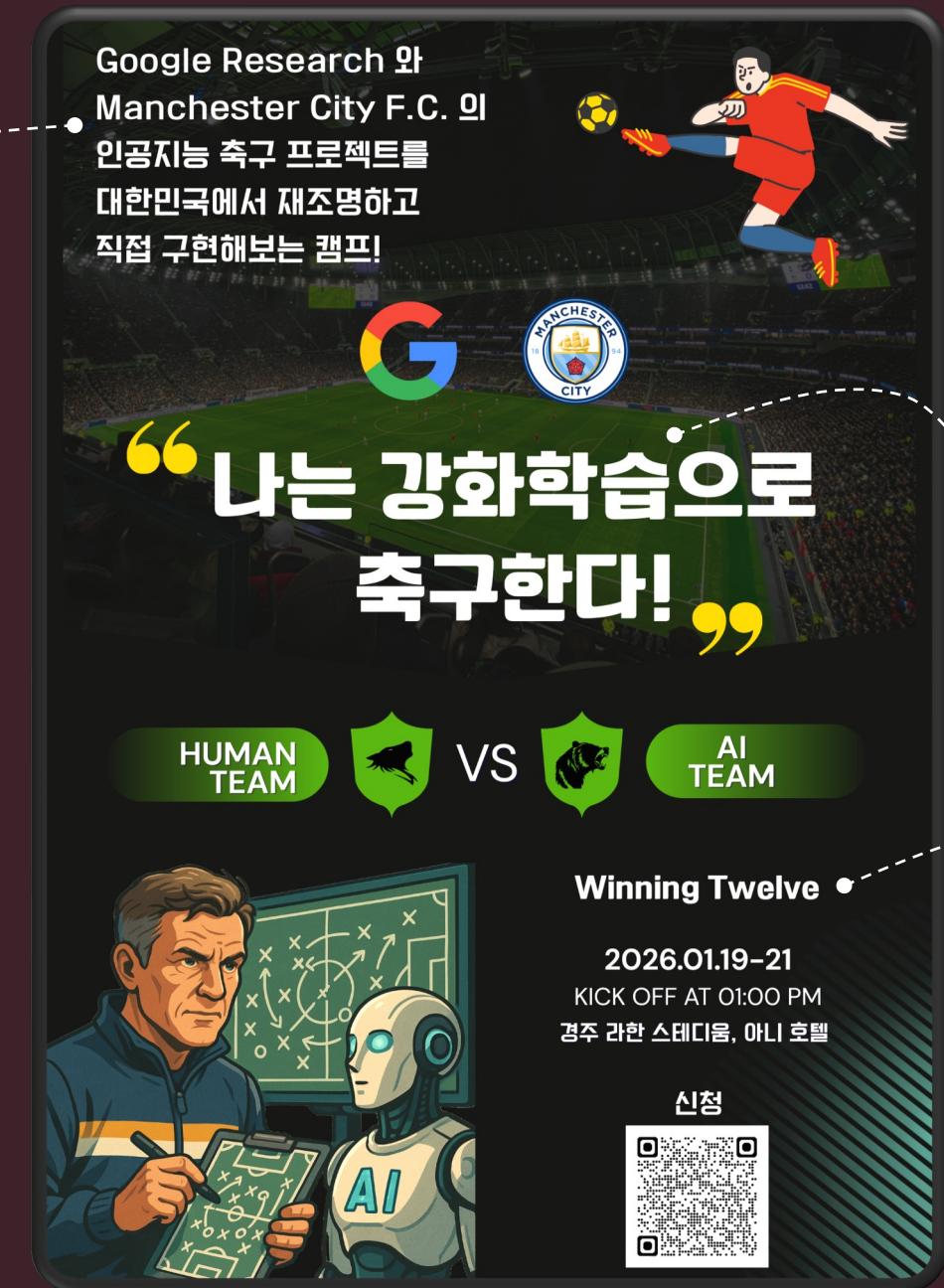
SESSION

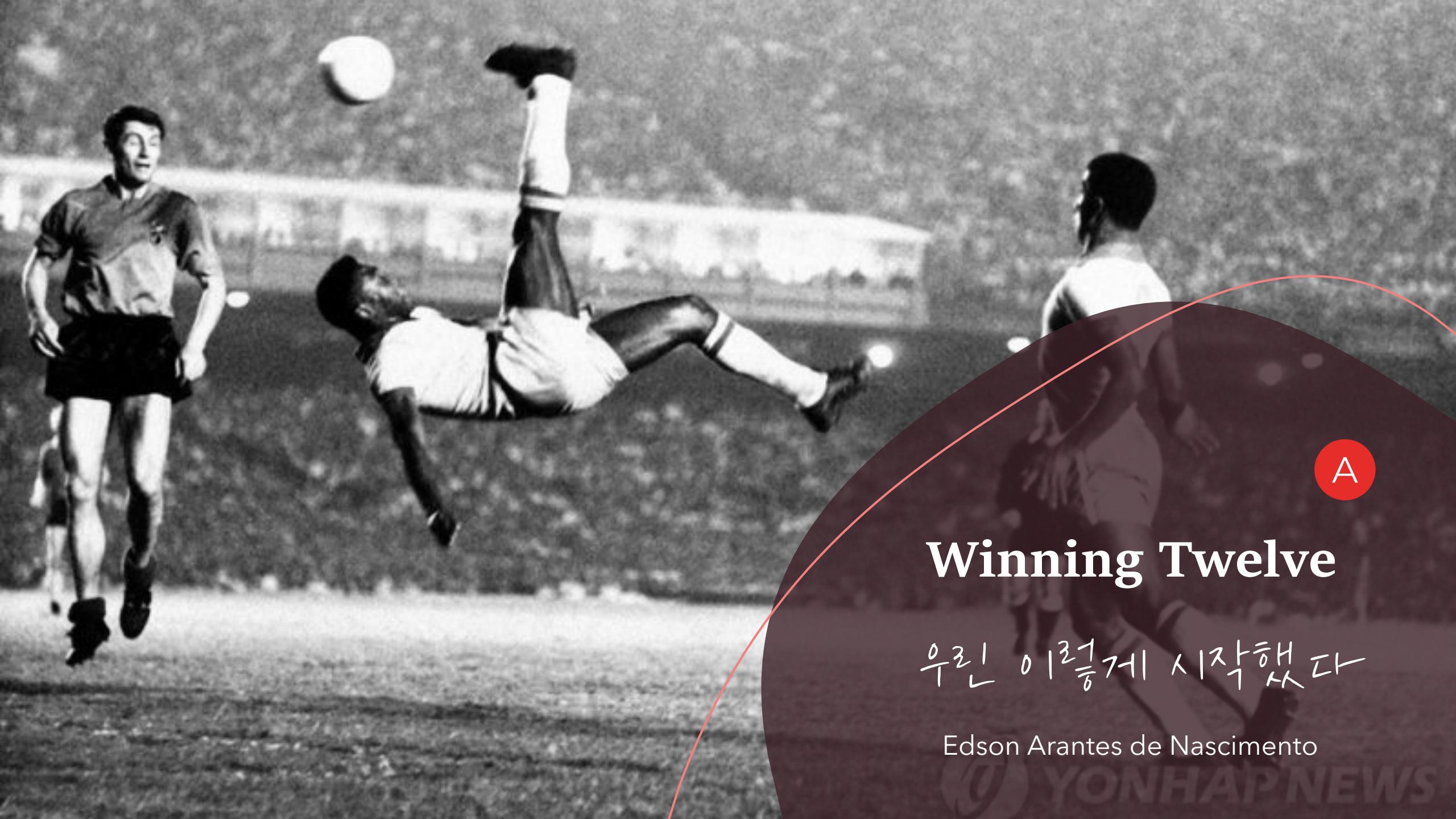
Copyright Notice

These slides are available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite the author as the source of the slides.

C

Google Research
Manchester City F.C.
재조명





A

Winning Twelve

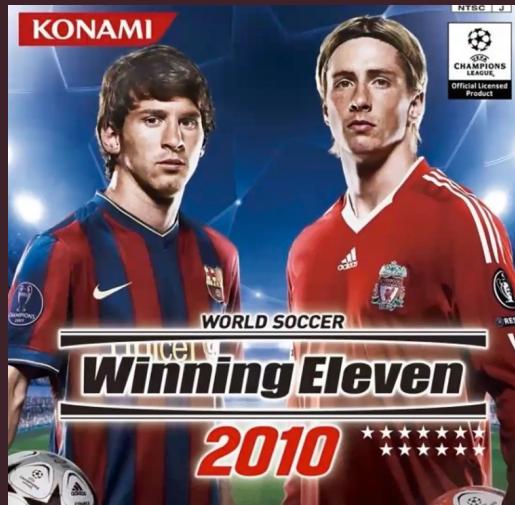
우린 0 1 $\frac{1}{2}$ 게 시작했다

Edson Arantes de Nascimento

YONHAP NEWS

A

Winning Twelve?



우린 0 1 $\frac{3}{8}$ 게 시작했다

- Edson Arantes de Nascimento
- Manchester United
- H-Milan

11 + 1



축구선수

데이터사이언티스트

Camp History

2023 Summer : 나는 데이터로 축구한다

2024 Summer : 나는 데이터로 축구한다

2025 Winter : 나는 데이터로 축구한다 (LLM/RAG 편)

2025 Summer : AI+X Ethiopia

2026 Winter : 나는 강화학습으로 축구한다



B

강화학습 / RL

Reinforcement Learning

- AlphaGo
- OpenAI Five (Dota 2)
- AlphaStar (Star Craft)
- OpenAI Gym
- ...

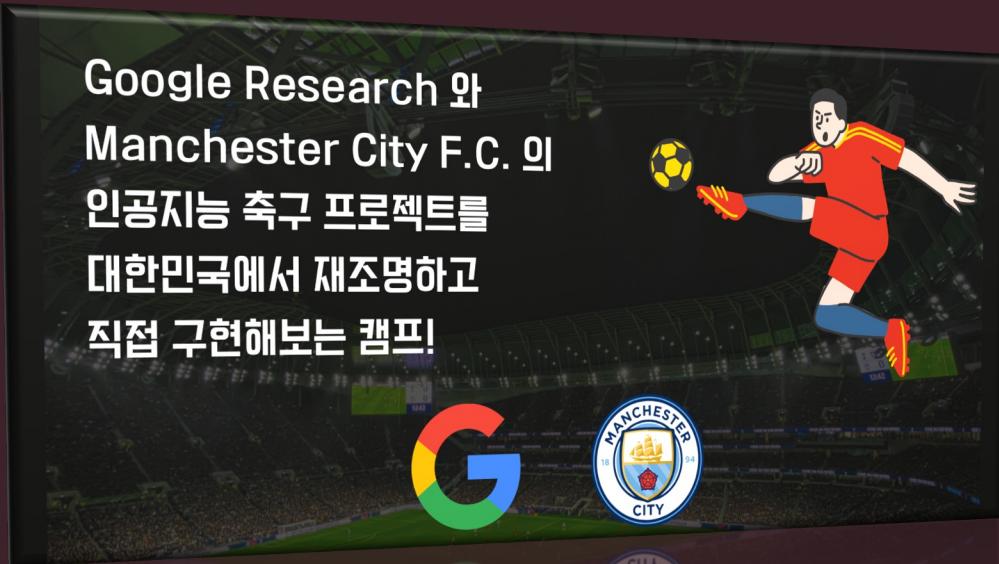


RL Applications

B

- Games
- Robot, Drone, Self-Driving Car
- Traffic Signal Control
- Route Optimization
- Fleet Management
- Intrusion Detection
- Recommendation Systems
- Algorithmic Trading
- Dynamic Pricing
- Drug Dosage Optimization
- Human-in-the-loop Learning (RLHF)
- ...

C



2025/12/03 부터

결혼, 계절학기, 인턴, 탈영(?), ...

Google Research Football with Manchester City F.C.

Leaderboard

This competition has completed. This leaderboard reflects the final standings.

Prize Winners

#	Team	Members	Score	Agents	Last	Solution
1	WeKick		1785.8	22	5y	View
2	SaltyFish		1597.5	15	5y	View
3	Raw Beast		1564.3	43	5y	View
4	zsp1197 zsp		1483.6	43	5y	View
5	TamakEri		1471.0	20	5y	View

수학과 알고리즘

Maximize $J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[G_0]$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) G_t \right] = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i) G_t^i$$

$$\theta \leftarrow \theta + \alpha \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i) G_t^i$$

$$\mathbf{a} \sim \pi_\theta(\cdot | \mathbf{s}_t)$$



- Markov Decision Process (MDP)
- Q-Learning
- Deep Q Network (DQN)
- Actor-Critic
- DDPG
- Proximal Policy Optimization (PPO)
- Impala CNN



Initialize policy π_θ and value function V_ϕ
for iter = 1..K:

$$\theta_{old} \leftarrow \theta$$

$$\phi_{old} \leftarrow \phi$$

Collect rollouts with old policy

$D = \{(s_{t,i}, \mathbf{a}_{t,i}, r_{t+1,i}, s_{t+1,i}, done_i)\}$
collected by $\pi_{\theta,old}$

Compute return and advantages

$$R_{t,i} = r_{t+1,i} + \gamma V_{old}(s_{t+1,i}) \times (\text{not } done_i)$$

$$\hat{A}_{t,i} = R_{t,i} - V_{old}(s_{t,i})$$

for epoch = 1..E:

for minibatch $B \subset D$:

Actor loss (clip objective)

$$r_{t,i} = \exp(\log \pi_\theta(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) - \log \pi_{\theta,old}(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}))$$

$$L_{clip} = \text{mean}_{i \in B} (\min(r_{t,i} \times \hat{A}_{t,i}, \text{clip}(r_{t,i}, 1-\varepsilon, 1+\varepsilon) \times \hat{A}_{t,i}))$$

Critic loss (MSE)

$$L_v = \text{mean}_{i \in B} (V_\phi(s_{t,i}) - R_{t,i})^2$$

Update

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} L_{clip}(\theta)$$

$$\phi \leftarrow \phi - \beta \nabla_{\phi} L_v(\phi)$$

Cost & Time



“

First, I tried a single V100 GPU, 12 vCPU cores, 32 GB of RAM, 128Gb regular disk. With this machine, I achieved about 17 SPS (Steps per second) and about 50-60 Game FPS (Frames per Second) for each environment. We'll run 16 environments in parallel, per environment, averaging ~280 FPS total. Estimated to take 50h and cost 108\$ for 50M steps.

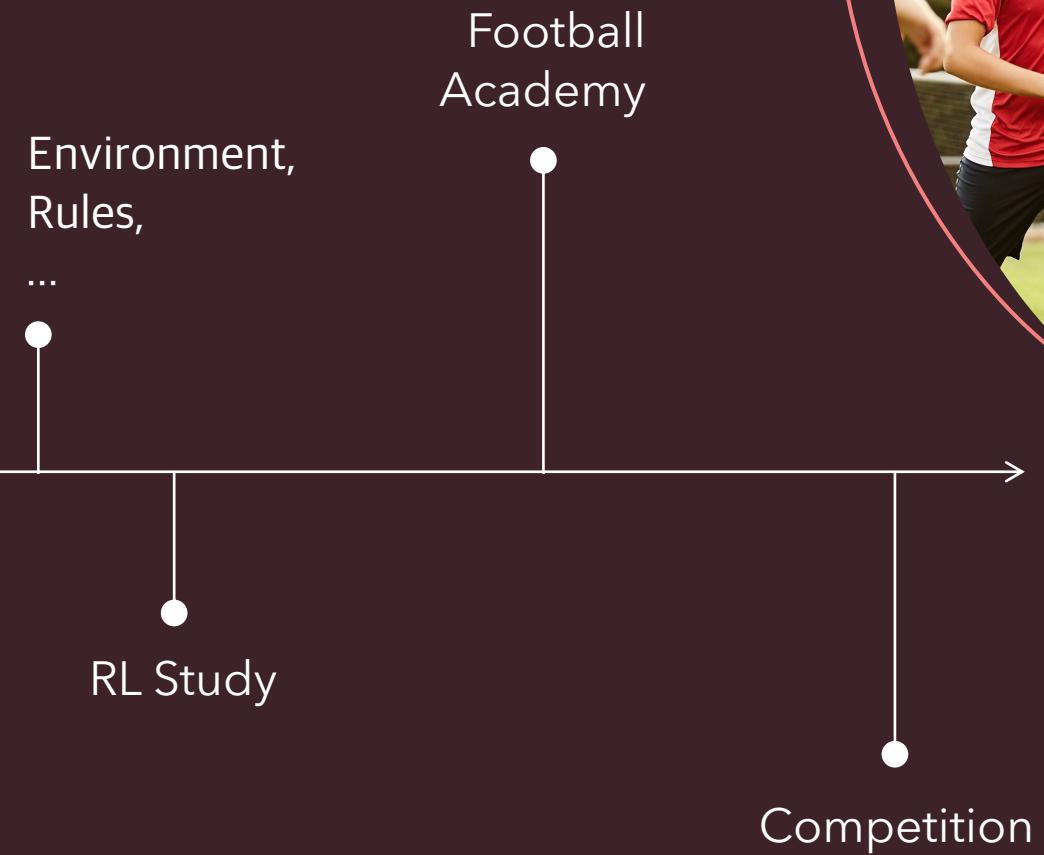
\$108

50 HOURS

시간이 없다.
돈도 없다.
쉽지 않다.



Agenda





00:01



FRQ 0



RBA 0

academy_empty_goal_close





00:01



FRQ



RBA 0

TURING

academy_empty_goal





00:01



FRQ



o



RBA



o

ARCHIMEDINHO

TURING

academy_run_to_score

00:00 FRQ 0 RBA 0

academy_run_to_score_with_keeper





00:01



FRQ



RBA



JIGSAW HD AVAILABLE SOON

Tenshu General

Tenshu General

BlauwPrint

BlauwPrint

JOHNSON



EINSTEIN

academy_pass_and_shoot_with_keeper





00:01



FRQ



RBA



JIGSAW HD AVAILABLE SOON

Tenshu General

Tenshu General

BlauwPrint

BlauwPrint

JOHNSON

EINSTEIN

academy_run_pass_and_shoot_with_keeper





00:01



FRQ 0



RBA 0

academy_3_vs_1_with_keeper



TURING

EINSTEIN





00:01



FRQ 0



RBA 0

academy_corner

BANNEKER

TURING

BlauwPrint**Indietopia®****Indietopia®**



00:01



FRQ o



RBA o

NEWTON

academy_single_goal_versus_lazy





00:01



FRQ



RBA



1_vs_1_easy





00:01



FRQ o



RBA o



5_vs_5



00:01 $\hat{\theta}$ FRQ 0 P RBA 0



11_vs_11_competition



Five Scenarios

1. academy_empty_goal_close
2. academy_empty_goal
3. academy_run_to_score
4. academy_run_to_score_with_keeper
5. 5_vs_5



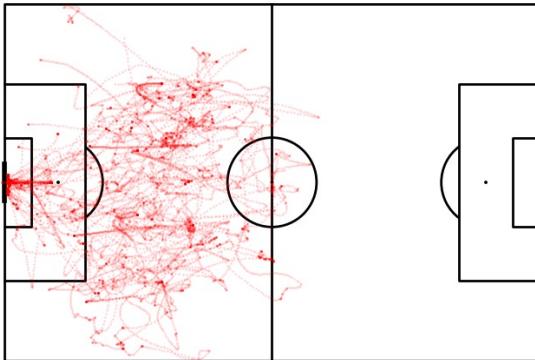
Winning Twelve World Cup

11-vs-11

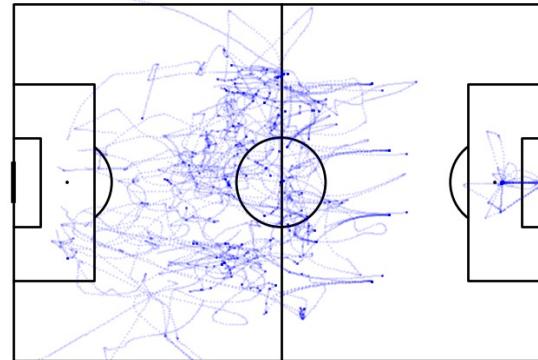
5-vs-5



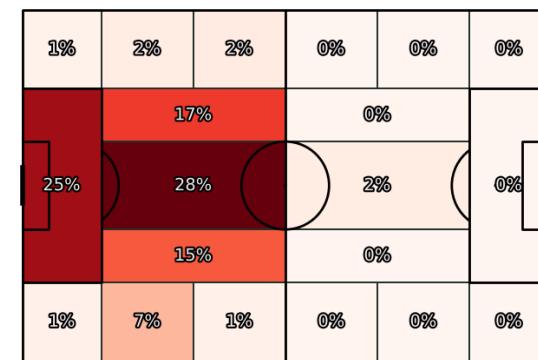
Your AI Team >>



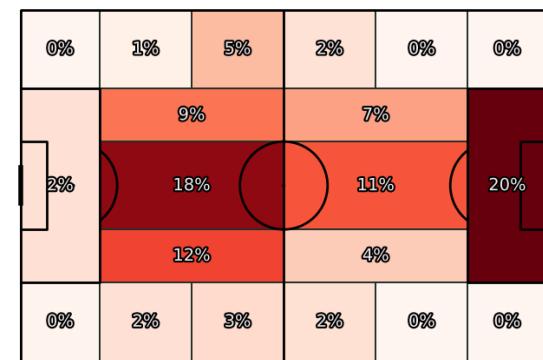
<< Google Research Team



Your AI Team >>



<< Google Research Team

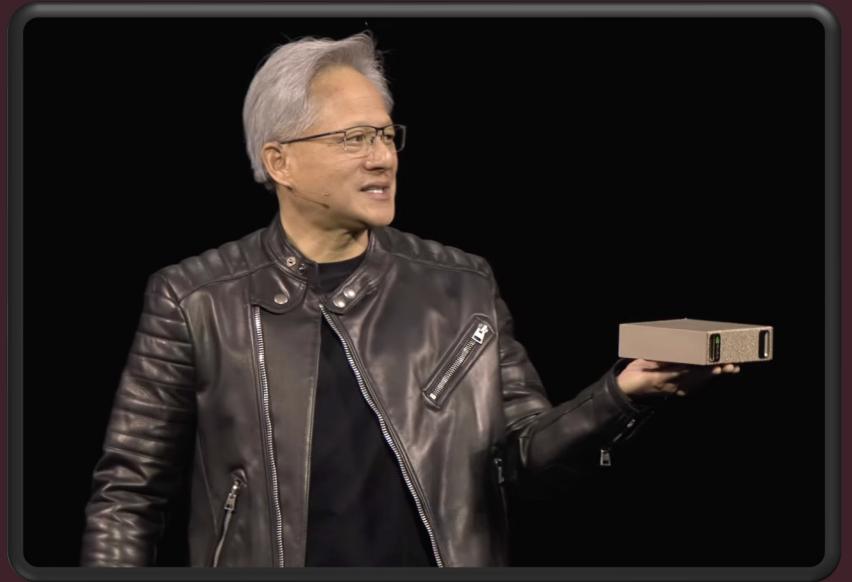


Leaderboard

- Ice breaking
- Run to Score with Keeper
- 시간엄수
- ...



NVIDIA DGX Spark @ GTC 2025 Keynote



Hand delivery to Elon Musk @ SpaceX

