**1** Collect resources

**2** Build a base

**3** Build units

**4** Defeat the opponent

**1** Complex Combinatorial Action Space

→ automatic theorem proving
→ drug design
→ industrial control problems
→ AGI
→ ...

**2** Multi-modal Observation Space

➔ robotics
➔ self-driving cars
➔ AGI
➔ ...

**3** Information "Poverty" and Hard Exploration

➔ natural sciences
➔ weather forecasting
➔ robotics
➔ AGI
➔ ...

**4**

Human "alignment"

→ self driving cars
→ home assistants
→ human enhancing AIs
→ AGI
→ ...

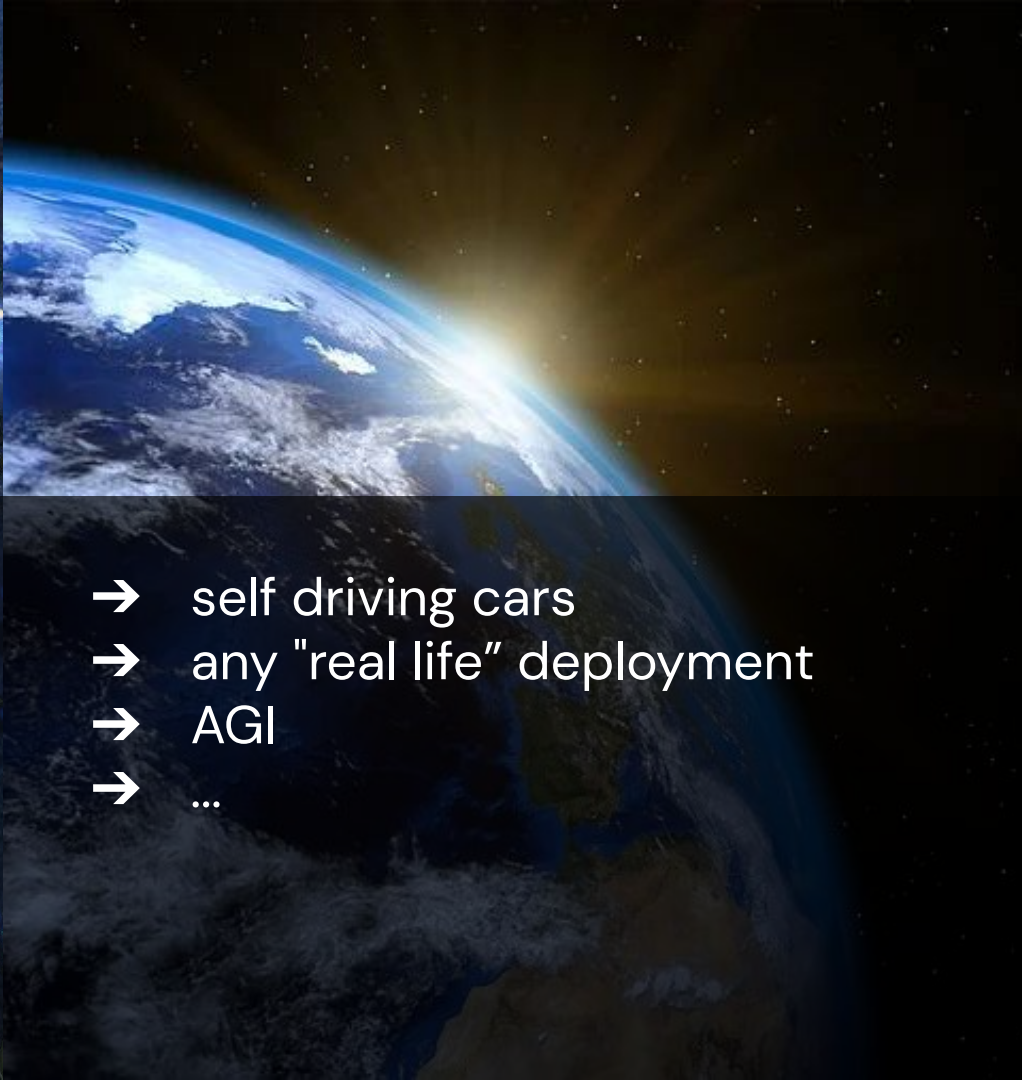**5** Multiple Interacting Agents

→ self driving cars
→ any "real life" deployment
→ AGI
→ ...

# YOU'VE BEEN PROMOTED!

IIIIIIIIIIII

## 1V1 GRANDMASTER

Your performance has qualified you for placement in a new league

OK

First season of SCII AI Ladder

Open-sourced SCII API

Sejong Uni comp all bots beaten by human

AlphaStar demonstration match, Jan 2019

AlphaStar GM level at the full game of SCII

ORTS

ORTS run first RTS games comp

First AIIDE comp

Top 3 bots beaten by human at AIIDE

ORTS

| 2003 | 2006 | 2009 | 2010 | 2011 | 2015 | 2016 | 2017 | 2018 | 2019 |

Academic advocacy of RTS games for AI

4th & final ORTS comp

Open-sourcing of BWAPI

IEEE CIG SC comp starts

SSCAIT comp starts

Announced SCII as research environment at BlizzCon

1V1
Go head to head in StarCraft II's competitive matchmaking

DEBUG

TERRAN

# DeepMind

## What's happening?

We're excited to announce that experimental versions of DeepMind's StarCraft II agent, AlphaStar, will soon play a small number of games on the competitive ladder as part of ongoing scientific research into artificial intelligence.

If you would like the chance to help DeepMind with its research by matching against AlphaStar, you can **opt in** by clicking the button below. If you opt-in and are matched against AlphaStar, DeepMind will use and may publish your match data and game replays in accordance with the terms below. Your username will not be published.  You can alter your opt-in selection at any time by using the "DeepMind opt-in" button on the 1v1 Versus menu.

For scientific test purposes, DeepMind will be benchmarking the system's performance by playing AlphaStar anonymously during a series of blind trial matches. This means the StarCraft community will not know which matches AlphaStar is playing, to help ensure all games are played under the same conditions. AlphaStar plays with built-in restrictions defined in consultation with pro players. A win or a loss against AlphaStar will affect your MMR as normal.

Thank you to everyone who has helped our work with DeepMind so far, and to all those who continue to support us as we push the boundaries of what's possible in StarCraft!

For more information on this work, review our FAQ here.

## Terms & Conditions

If you are matched against AlphaStar, DeepMind Technologies Limited (a company organised under the laws of England and Wales) will use the games, game replays and game data created to conduct research on the development of machine learning, which may include publication of some replays

OPT-IN                 OPT-OUT

RANKED          UNRANKED          MAPS

DeepMind

# Interface

## Human alignment

Actions limit ~22 per 5 s — Requested delay ~200 ms

**Action**

Move
Attack
Build

What? → Who? → Where? → When next action?

**Own units**
Camera vision | Outside camera

**Opponents units**
Camera vision | ? Outside camera

**Minimap**

Move

# Reduced APM - Pro tested and approved

● AlphaStar  ● Battle.net opponents



APM

APM

APM

Zerg has higher APMs due to
repeat actions, such as morphing
& spawning

DeepMind

# Supervised learning

**Hard exploration**
**Information poverty**

$a_t$

KL

$\pi_t^{SL}$

$o_t$

Human

t

Z

Supervised 936

No human data 149

0    600    1200    1800    2400

Test Elo

13/14    17/22    18/22    17/22    19/22    19/22    21/30

Even AlphaStar Supervised is not a single "strategy". It is a (controllable!) collection of dozens of thousands of strategies

DeepMind
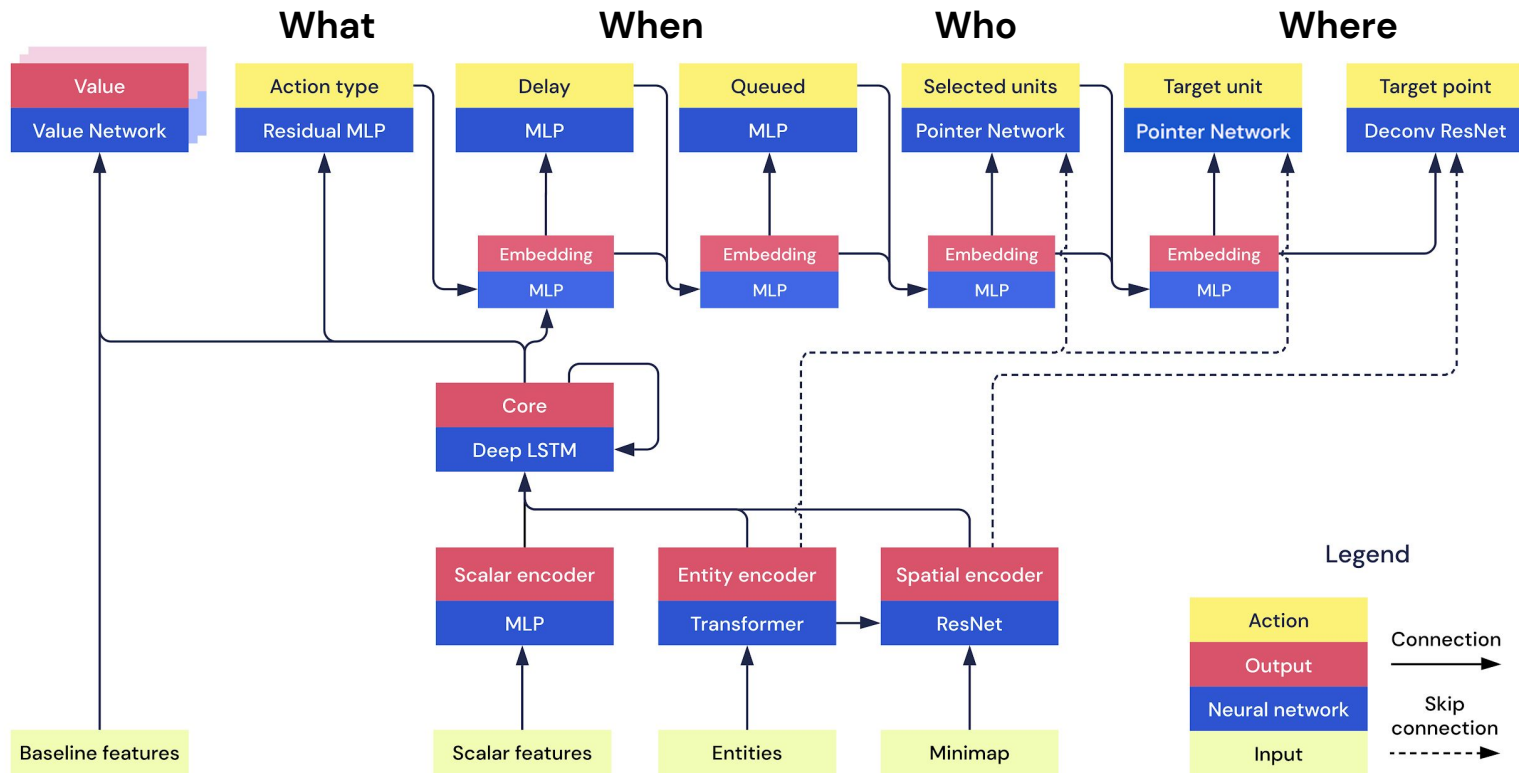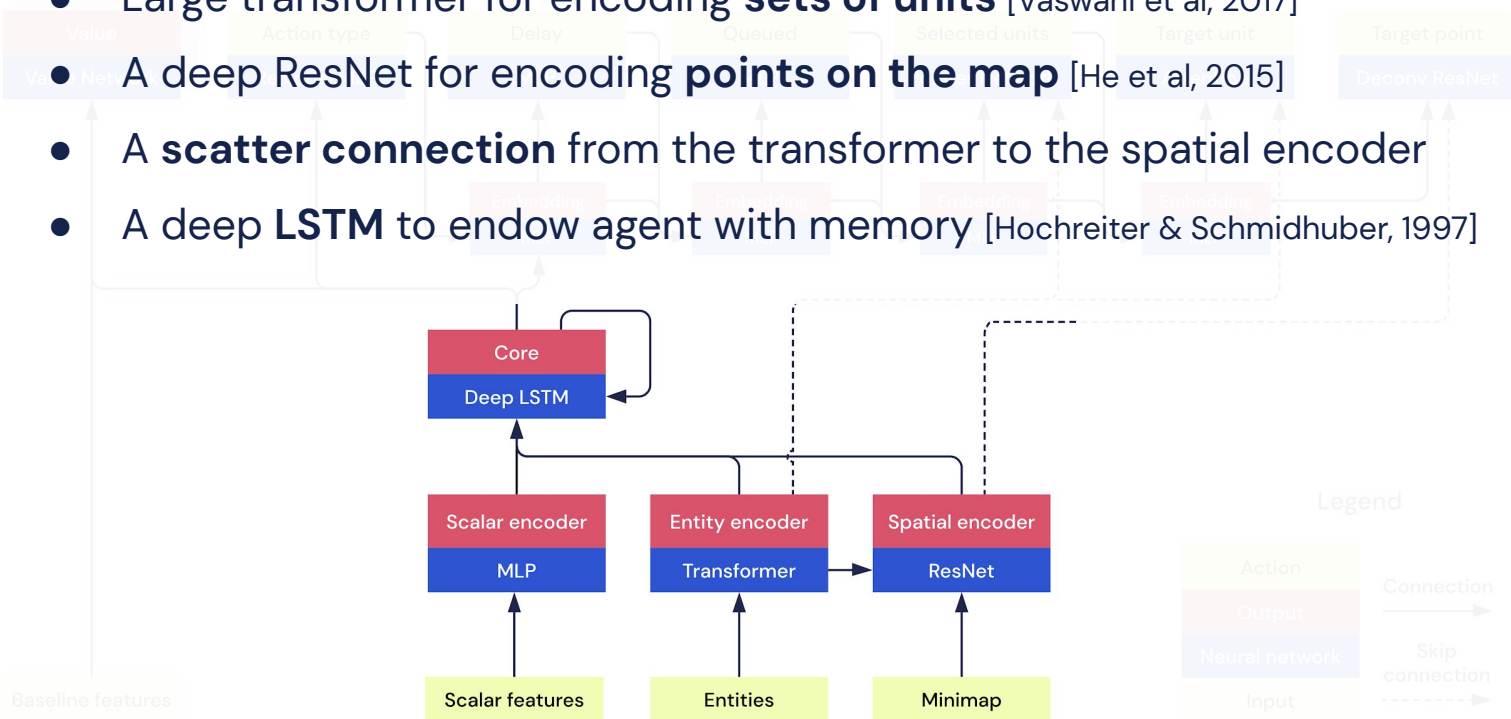
# Architecture
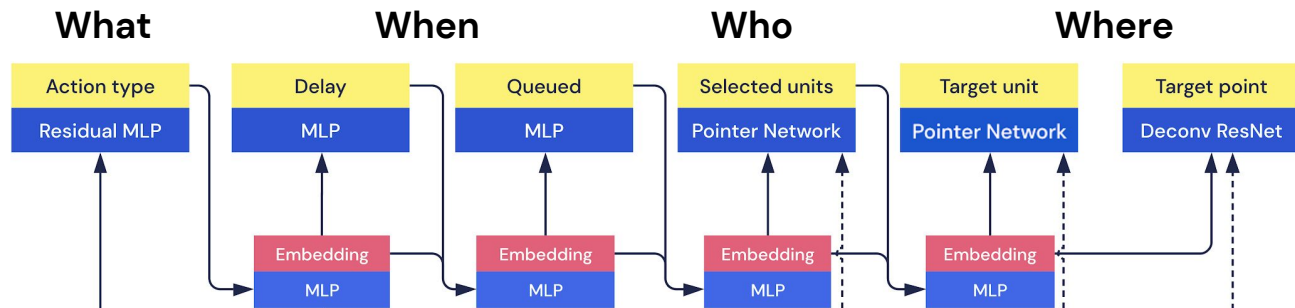## Combinatorial Action Space
## Multi-modal Observations

# Observation encoders and LSTM

- Large transformer for encoding **sets of units** [Vaswani et al, 2017]
- A deep ResNet for encoding **points on the map** [He et al, 2015]
- A **scatter connection** from the transformer to the spatial encoder
- A deep **LSTM** to endow agent with memory [Hochreiter & Schmidhuber, 1997]

# Autoregressive action head

| What | When | | Who | | Where |
|------|------|------|-----|------|-------|

**What**
- Action type / Residual MLP

**When**
- Delay / MLP
- Queued / MLP

**Who**
- Selected units / Pointer Network
- Target unit / Pointer Network

**Where**
- Target point / Deconv ResNet

Embedding / MLP (repeated across When, Who, Where columns)

- Fully autoregressive action head with 7 sub–heads: $p(x) = \prod_{i=1}^{n} p(x_i | x_1, \ldots\ldots, x_{i-1})$

- Four scalar heads: **action type**, **action delay**, **action repeat, modifier key**

- A recurrent pointer network to select a **set of units**  [Vinyals, Fortunato & Jaitly, 2015]

- A simple pointer network to select **single units**

- A ResNet decoder to select **points on the map**

# Story time:
## *Why is our Zerg agent so bad compared to other races?*

- Other races are just better initially, they dominate the Zerg and it cannot flourish

- We are missing some of the Zerg-specific observations/parts of environment

# Story time:
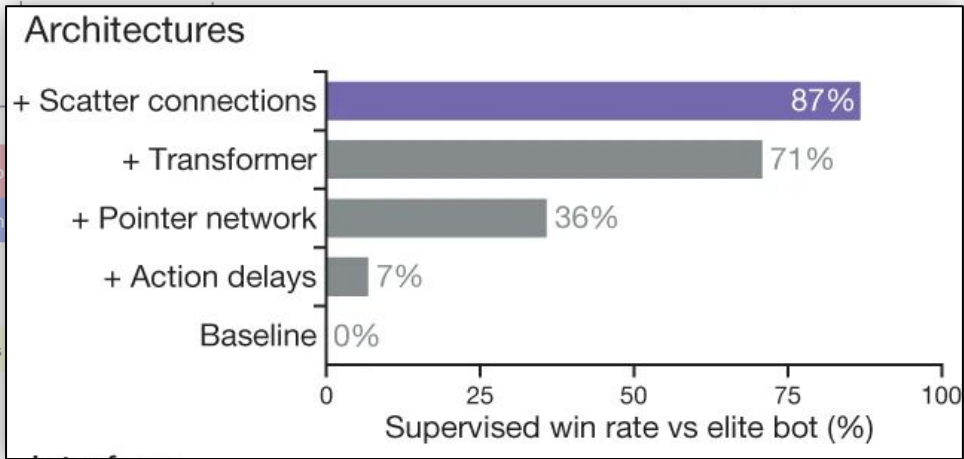## *Why is our Zerg agent so bad compared to other races?*
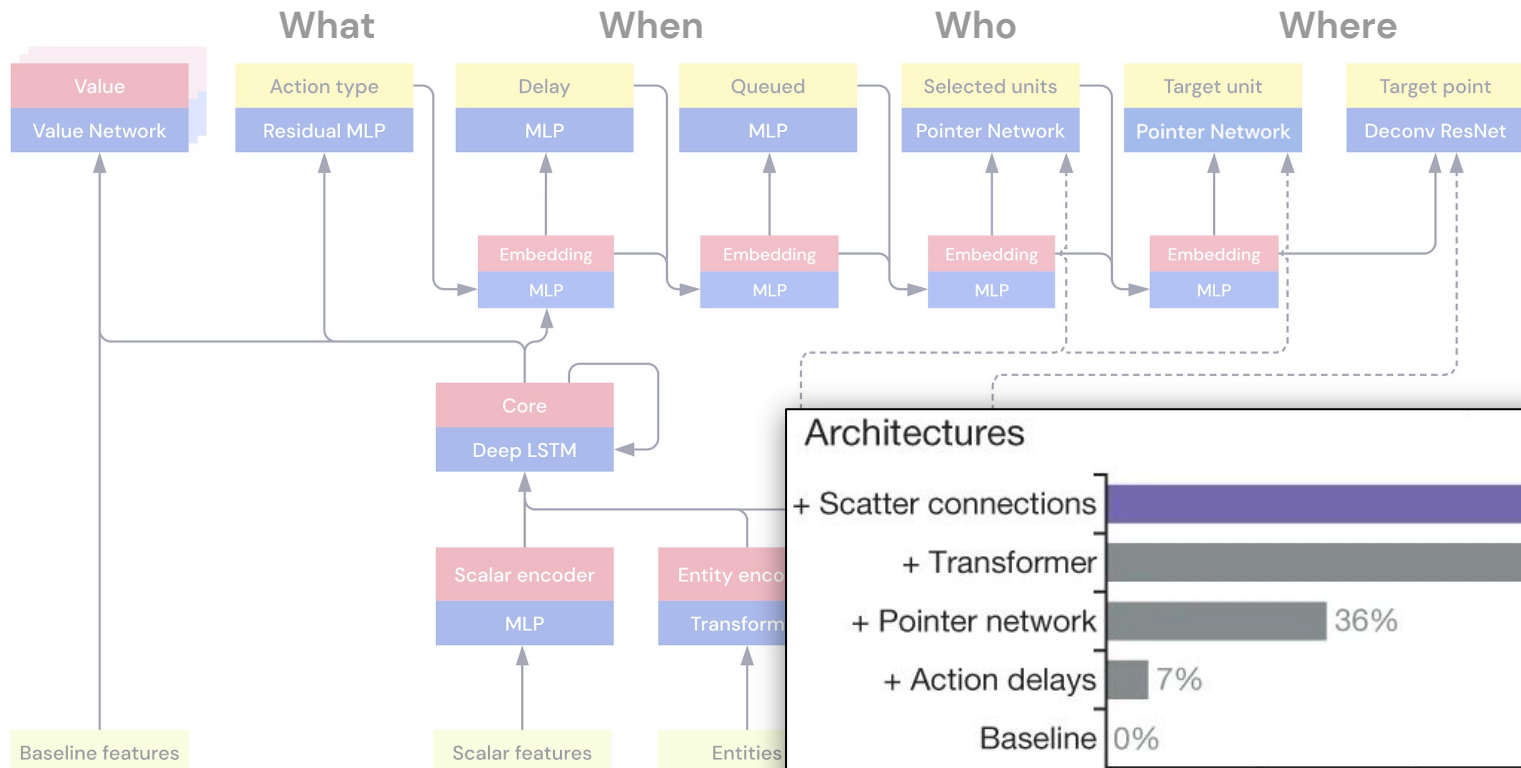
- Other races are just better initially, they dominate the Zerg and it cannot flourish

- We are missing some of the Zerg-specific observations/parts of environment

- It's because our Protoss agent is not good enough, and there is an architectural trick missing!

# Story time:
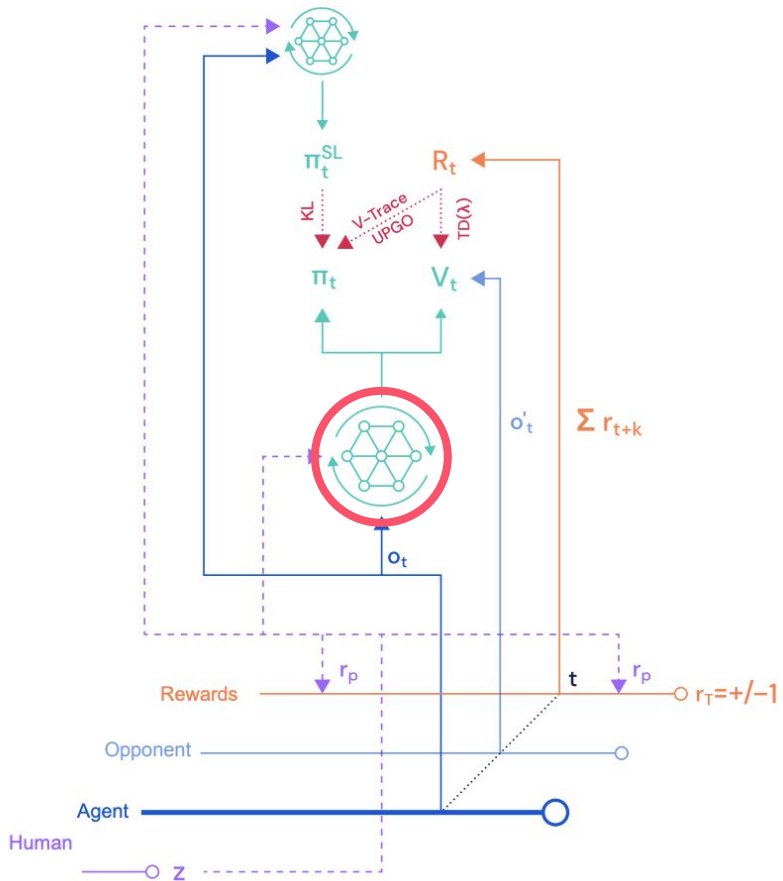## *Why is our Zerg agent so bad compared to other races?*

**What**

| Value | Action type |
|---|---|
| Value Network | Residual MLP |

**When**

| Delay | Queued |
|---|---|
| MLP | MLP |

**Who**

| Selected units | Target unit |
|---|---|
| Pointer Network | Pointer Network |

**Where**

| Target point |
|---|
| Deconv ResNet |

| Embedding | Embedding | Embedding | Embedding |
|---|---|---|---|
| MLP | MLP | MLP | MLP |

| Core |
|---|
| Deep LSTM |

| Scalar encoder | Entity enco... |
|---|---|
| MLP | Transform... |

Baseline features

Scalar features

Entities

### Architectures

+ Scatter connections — 87%
+ Transformer — 71%
+ Pointer network — 36%
+ Action delays — 7%
Baseline — 0%

Supervised win rate vs elite bot (%)
0    25    50    75    100

## Reinforcement Learning

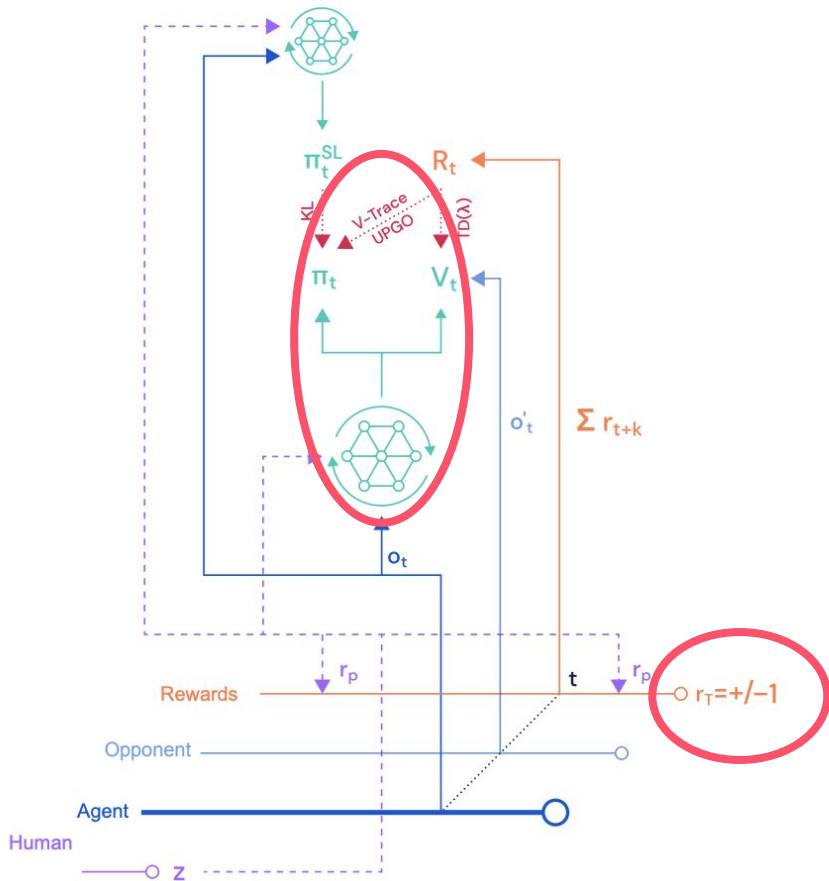$\pi^{SL}$

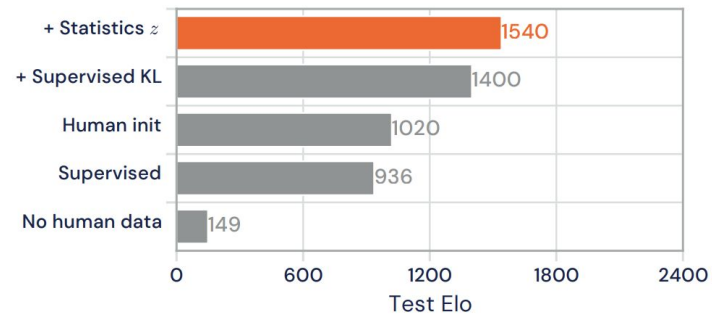### E  Human data usage

Test Elo

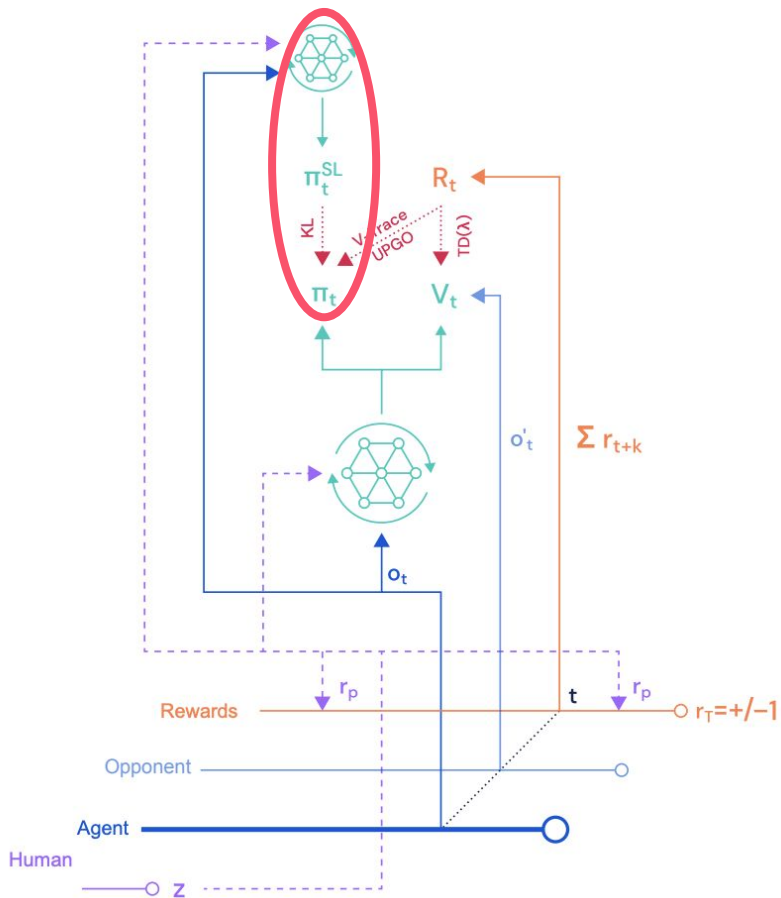## Reinforcement Learning

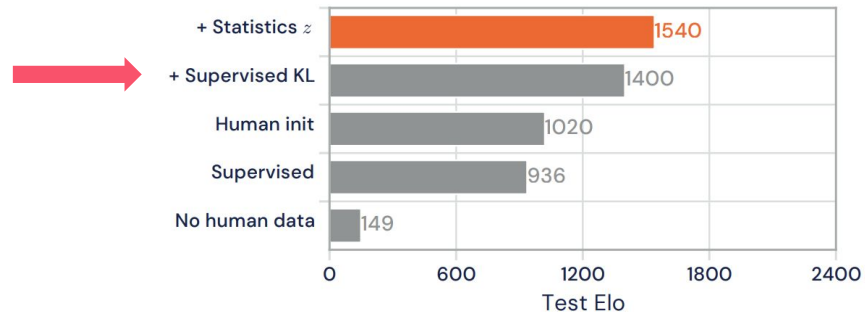$$\pi_0 = \pi^{SL}$$

E   Human data usage

## Reinforcement Learning

$$\pi_0 = \pi^{SL}$$
$$KL(\pi, \pi^{SL})$$

### E   Human data usage
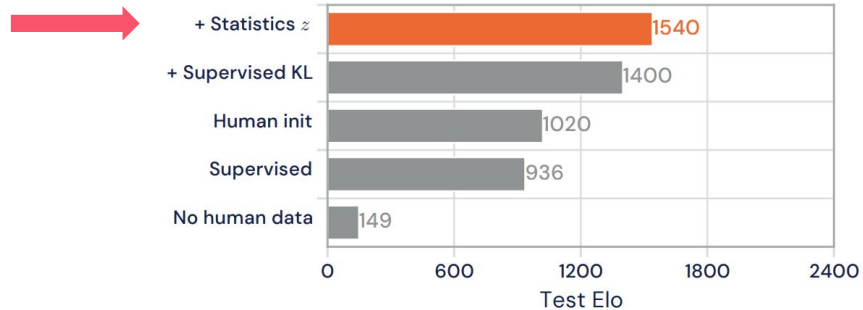
# Reinforcement Learning



## E  Human data usage

Do not start thinking about multi-agent dynamics research until you have a fully working, robust "best response" setup.
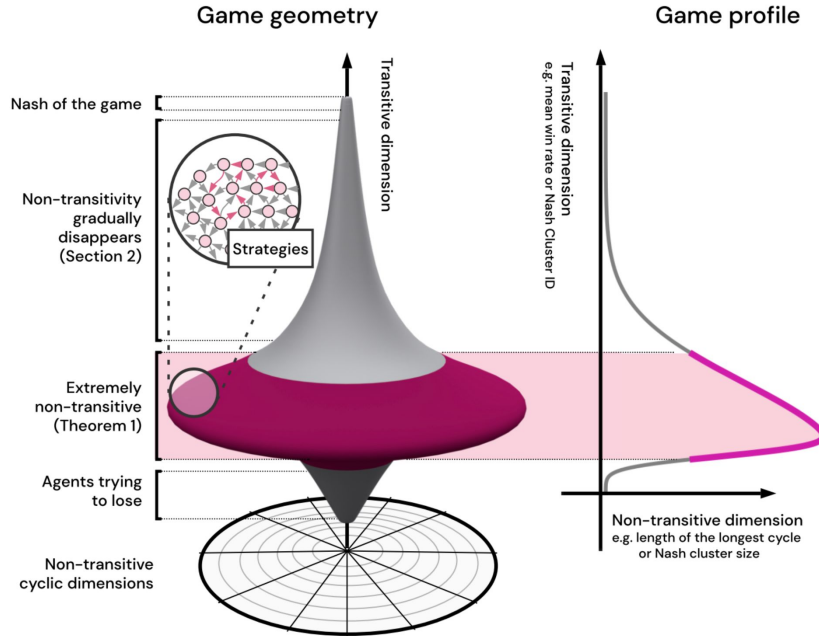
DeepMind

# Multi-agent Learning

**Multiple Interacting Agents**
**Hard exploration**
**Information poverty**

# Real world games look like spinning tops



Game geometry

Game profile

# Payoff matrix analysis

# Rock



## Scissors

## Paper

**Rock, paper, scissors**

StarCraft II players can create a variety of 'units', which have balanced strengths and weaknesses, similar to the game rock, paper, scissors

## Self-play



Starts playing a mixture of units

Specialises

Specialisation becomes narrower

- **V** Void ray
- **S** Stalker
- **I** Immortal
- ← Strong against
- Exploiter units
- Agent units
- — Distribution of units that agent can beat with its current strategy

Rock

Scissors

Paper

**Rock, paper, scissors**

StarCraft II players can create a variety of 'units', which have balanced strengths and weaknesses, similar to the game rock, paper, scissors

**Self-play**

Starts playing a mixture of units

Specialises

Specialisation becomes narrower

**With exploiters**

Starts playing a mixture of units

Exploiters highlight weakness

Pursues more robust strategies

V  Void ray
s  Stalker
i  Immortal
←  Strong against
●  Exploiter units
●  Agent units
—  Distribution of units that agent can beat with its current strategy

The main agent fails to defend itself.

# The League Training



→ Build strong, robust agents

→ Expose main agents weaknesses

→ Expose weaknesses of entire League

# The League Training

# The League Training

# The League Training

## A  League composition



Test Elo

- + League Exploiters: 1824
- + Main Exploiters: 1693
- Main Agents: 1540

## B  League composition



Relative Population Performance

- + League Exploiters: 62%
- + Main Exploiters: 35%
- Main Agents: 6%

# Who to train against?

$$\forall_i \ \mathrm{P}(\pi_\theta \text{ winning against } \pi_i) > 0.5$$

# Who to train against?

$$\forall_i \ \mathrm{P}(\pi_\theta \text{ winning against } \pi_i) > 0.5$$

→     FSP     $\mathrm{U}\big(\{\pi_i\}_{i=1}^{N}\big)$

# Who to train against?

$$\forall_i \; \mathrm{P}(\pi_\theta \text{ winning against } \pi_i) > 0.5$$

→ FSP $\quad \mathrm{U}\big(\{\pi_i\}_{i=1}^{N}\big)$

→ **P**FSP $\quad \mathrm{P}(\text{playing against } \pi_j) = \dfrac{f(\mathrm{P}(\pi_\theta \text{ winning } \pi_j))}{\sum_i f(\mathrm{P}(\pi_\theta \text{ winning } \pi_i))} \qquad f : [0, 1] \to \mathbb{R}_+$

# Who to train against?

$$\forall_i \ \mathrm{P}(\pi_\theta \ \text{winning against} \ \pi_i) > 0.5$$

→ FSP $\quad \mathrm{U}(\{\pi_i\}_{i=1}^N)$

→ **P**FSP $\quad \mathrm{P}(\text{playing against} \ \pi_j) = \dfrac{f(\mathrm{P}(\pi_\theta \ \text{winning} \ \pi_j))}{\sum_i f(\mathrm{P}(\pi_\theta \ \text{winning} \ \pi_i))}$ $\qquad f : [0, 1] \to \mathbb{R}_+$

$$f_{\mathrm{hard}}(x) = (1 - x)^p$$

You **never** play against opponents that you **dominate**.

You **focus** on **beating everyone** rather than average win rate.

When a **rare**, but **strong**, opponent appears – it is being **focused** on.

# Who to train against?

$$\forall_i \ \mathrm{P}(\pi_\theta \text{ winning against } \pi_i) > 0.5$$

→ FSP $\qquad \mathrm{U}(\{\pi_i\}_{i=1}^N)$

→ **P**FSP $\qquad \mathrm{P}(\text{playing against } \pi_j) = \dfrac{f(\mathrm{P}(\pi_\theta \text{ winning } \pi_j))}{\sum_i f(\mathrm{P}(\pi_\theta \text{ winning } \pi_i))} \qquad\qquad f : [0,1] \to \mathbb{R}_+$

$$f_{\mathrm{hard}}(x) = (1-x)^p$$

You **never** play against opponents that you **dominate**.

You **focus** on **beating everyone** rather than average win rate.

When a **rare**, but **strong**, opponent appears – it is being **focused** on.

$$f_{\mathrm{var}}(x) = (1-x)x$$

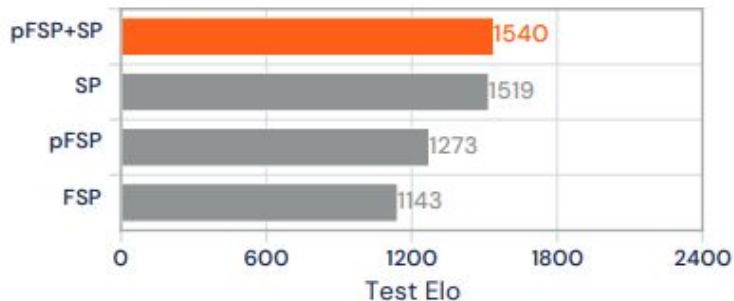You always pick **opponents** at your **own level**.

Creates a natural auto **curriculum**.

A **black-box** version of **TD-error prioritisation**.

# Who to train against?



C  Multi-agent learning

| | Test Elo |
|---|---|
| pFSP+SP | 1540 |
| SP | 1519 |
| pFSP | 1273 |
| FSP | 1143 |

D  Multi-agent learning

| | Min win-rate vs. past |
|---|---|
| pFSP+SP | 71% |
| SP | 46% |
| pFSP | 70% |
| FSP | 69% |

# Path matters

→ There are often infinitely many solutions for "best response to a fixed set of opponents" problem

→ "Greedy" decisions on the way identify which one we will end up with

→ They might differ dramatically with respect to their transitive strength and properties.

→ "Hard opponnents" can prefer policies of low transitive strength, but converges fast and to diverse policies.

$$f_{\text{hard}}(x) = (1 - x)^p$$

→ "Variance" produces more "standard" strategies, but converges much slower (and somewhat deterministically).

$$f_{\text{var}}(x) = (1 - x)x$$

# Path matters



G APM limits

| | Test Elo |
|---|---|
| No APM limit | 1392 |
| 200% APM limit | 1411 |
| 100% APM limit | 1540 |
| 50% APM limit | 1536 |
| 25% APM limit | 1419 |
| 10% APM limit | 1145 |
| 0% APM limit | 0 |



Control

Complex counter

Trivial counter

Strategy A    Strategy B

# Mutli-agent Deep Reinforcemenet Learning

$$!=$$

**Multi-agent** +
**Reinforcement learning** +
**Deep learning**

# It is a "new field"

- Choice of correct opponents is not just guided towards convergence to the Nash, but also takes into consideration dynamics of Deep RL
- RL needs to be curated towards specifics of Multi-agent, e.g. rapid changes of targets, non-stationarity
- Exploration is not just an RL issue, with multi-agent algorithms we can guide the weak exploration strategy to shine in a complex problem
- Architectures can create entire new levels of multi-agency
- Architectures, and improvements that are the best in simplified setups are not the ones that shine in the long term - speed of convergence is a wrong thing to optimise!
- Even the game interface shapes the dynamics of RL, and multi-agent!

|   | | |
|---|---|---|
| **1** | Complex Combinatorial Action Space | ➔ human-like constraints |
| **2** | Multi-modal Observation Space | ➔ architecture |
| **3** | Information poverty and Hard Exploration | ➔ new RL objectives |
| **4** | Human "Alignment" | ➔ "human exploration" |
| **5** | Multiple Interacting Agents | ➔ AlphaStar League |
|   | | ➔ A lot of hard teamwork! |

# Questions?