

O. Vinyals  
(DeepMind)  
- AlphaStar -

S. Whittaker  
(UoOxford)  
- Bayes-Adaptive  
RL -

# New RIPS 2019

## - Day 6 -

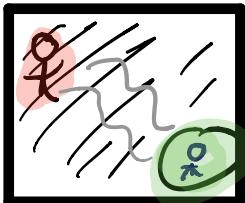
## WORKSHOPS II

David Ha  
(Google Brain)  
- Generative -

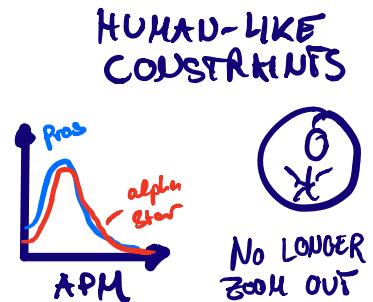
Oriol Vinyals (DeepMind): AlphaStar - Grandmaster  
Level in StarCraft II using Multi-agent RL

# STARCRAFT II → A NEW CHALLENGE

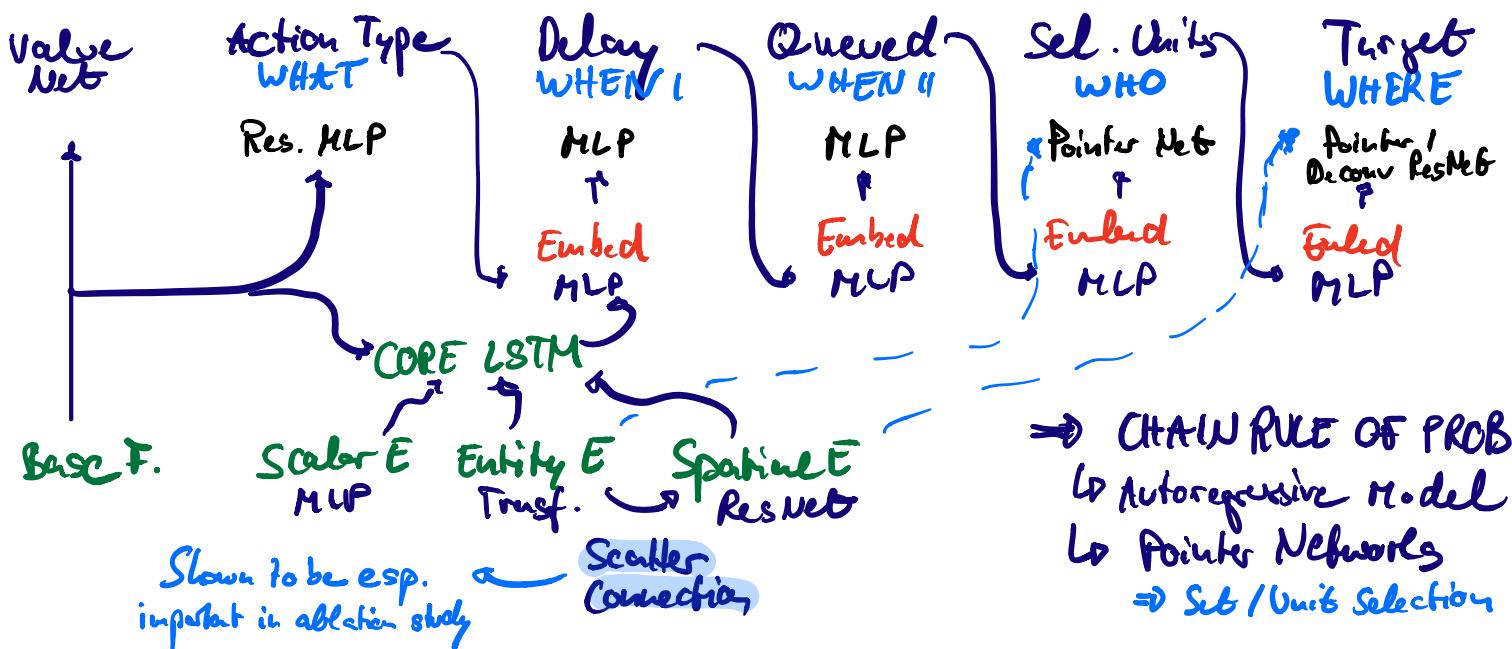
- ① Collect Resources
  - ② Build base
  - ③ Build units
  - ④ Defeat agents



- \* Partially Obs.
  - \* Multigener
  - \* 1000s of hours
  - \* Range  $\approx 10^{20}$



ACTIONS + ARCHITECTURE → SEQ. MODELLING PARALLELS



**EXPLORATION** → BEHAVIORAL CLONING + PRIOR REGULARIZ.

- Long timescales very challenging
  - Naive Worker Push  $\Rightarrow$  local Optimum

$\Rightarrow$  Instead: IMITATION LEARNING  
VI & BEHAVIORAL CLONING

  - Ⓐ  $\min KL(\pi^{\text{expl}} || \pi^{\alpha^*})$
  - \* Embedding Context  $Z$ 
    - ↳ Only during training

- RL Training  $\oplus$  KL prior on expt behavior traces
  - TD( $\lambda$ ) version: UPBO // V-trace  
↳ Off vs. On-Policy

# ROBUSTNESS → LEAGUE

- 'The League': Stere Chkps.
  - Enforce diversity!
  - Hold-Out agent validation

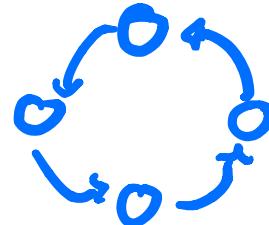
# Shimon Whiteson (MoOxford) : Bayes-Adaptive DRL via Meta-L.

EXPLORATION  $\rightarrow$  NAIVE  $\Rightarrow$  NO CONDITIONING OF BEHAVIOUR ON UNCERTAINTY EST.

Systematic Exploration :  $R_{\max} \rightarrow$  Family PTC methods  $\Rightarrow$  Theoretical Guarantees

Poor init rewards  
Over-Exploration?!  
Tabular Method

Pseudo-Count  
 $\Leftrightarrow$  Density  
Method



REAL

- Rewards Real/Cont
- Balance Crucial
- Regret Min.

SIMULATED

- Save env
- Sample-efficient
- Best consider.

$\Rightarrow$  CLAIM : There are no good methods for real exploration at the moment!

Principled Use of  
Prior Structure (Model)

Reward-traced  
Dirichlet Metacount

## BAYESIAN RL

① Learning  $\Rightarrow$  Task Uncertainty

② Bayes-Adaptive  $\Rightarrow$  Hyper-State

$$S^+ = S \times \mathbb{B} \text{--- Belief Space}$$

$\hookrightarrow$  Explore Only if it decreases task uncertainty + reward increase

## CHALLENGES

③ T, R representation

④ Intractable Planning

Instead:  
Posterior Sampling

⑤ Intractable Inference

Plan only in small MDP

$\Rightarrow$  No longer Bayes Optimal

variBAD  $\Rightarrow$  Meta-L. Framework

$$\pi_i = p(\mu_i) \rightarrow \text{Content } \alpha_i$$

$$\text{CD Infer } p(\mu_i; \mathcal{H}_i)$$

$\Rightarrow$  LEARN TO APPROX. B-OPT. BEH.

$\hookrightarrow$  Work by Luisa Zintgraf

ENCODER  $\rightarrow$  Decoder  $p_\theta$   
 $q_\phi \sim m$   $\xrightarrow{\quad}$  policy  $\pi_\phi$

$\hookrightarrow$  Bayes-Adaptive Perspective on Meta-RL

$\hookrightarrow$  Gridworld Validation  $\Rightarrow$  Decoder Interpretation of Learn. Dynamics

$\hookrightarrow$  PEARL: Rakelly et al., 18' - Off-Policy

$\hookrightarrow$  Info Bottleneck Interpretation!

$\hookrightarrow$  Half-Cheetah Val.: Quickly deduce task!

# David Ha (Google Brain): Innate Bodies, Brains & Minds

$\alpha$  DEAD FISH = BODIES AS EVOLVED PRIORS  
 $\alpha$  ↓

R2 for Improv. Agent Design  
 (Ha et al., 18')

DIRECTLY OPTIMIZE CONSTRAINTS OF AGENTS

Anthony Zador: 'Critique of Pure Learning [...] → Nus. N'

→ EVOLUTIONARY PERSPECTIVE ON INDUCTIVE BIASES

→ BUILDING BLOCKS



( ) COMPLEX ENV  $\leftrightarrow$  ROBUST POLICY?  
 ⇒ Path of least Resistance

Deep Image Prior  
 ulyanov et al. 18'

PASSIVE WALKERS  
 McGeer 90'  
 Collins et al. 01'

Schmidhuber - LSTMs  
 ↳ bias towards seg. 2.

## FINDING GOOD PRIORS FOR DRL → WEIGHTS (Gaior & Ha 19')

- ① Search for architectures by de-emphasizing weights  $\Rightarrow$  RANDOMIZE
- ② NEAT-based evolutionary algo  $\Rightarrow$  Fast tuning of single tuned weight
  - ↳ Bipedal Walker: „Learn“ to ignore parts of inputs
  - ↳ MNIST: Ensemble-based classification

① INIT  $\rightarrow$  ② EVALUATE  $\rightarrow$  ③ RANK  $\rightarrow$  ④ TRY

$\Rightarrow$  AUTOMATED DISCOVERY - Focus on NEW BUILDING BLOCKS

OPERATORS  
 INSERT NODE  
 ADD CONNECT.  
 CHANGE ACTIV.

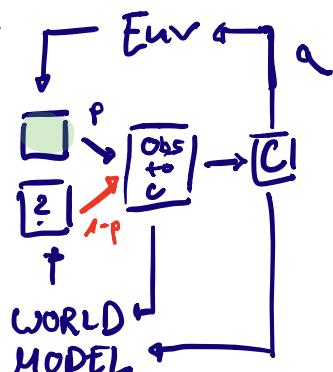
## PREDICTING THE FUTURE PRIOR $\rightarrow$ OBSERVATIONAL DROPOUT

Ha & Schmidhuber 18'  
 Rec. World Models

Scaling not goal  
 by itself

$E \dashrightarrow E^*$   
 TRANSFER  
 "WAKE ON, WAKE OFF"

- ① Freeman et al. 19':  
 $1-p\%$  access  $\Rightarrow$  Blindfold "agent"  
 ↳ learned model only respects  
 "important" aspects of env!  
 $\Rightarrow$  NOT PHYSICS ENGINE



'A CREATURE DIDN'T THINK IN ORDER TO MOVE; IT JUST MOVED, AND BY MOVING IT DISCOVERED THE WORLD THAT THEN FORMED THE CONTENT OF ITS THOUGHT.' - Andy Clark 08'