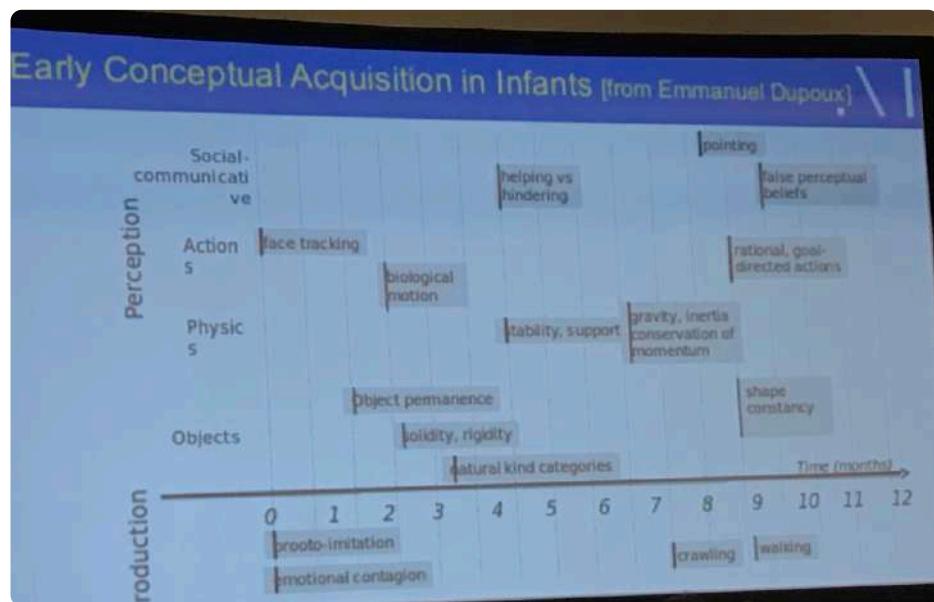


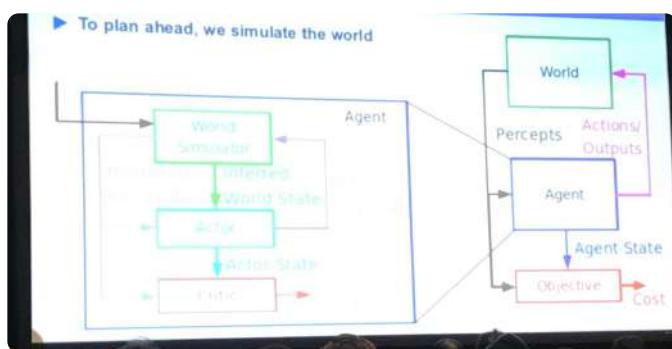
NIPS - Day 5 Workshops I - Deep Reinforcement Learning

① Y. LeCun - Learning Predictive Models

- Sample complexity problem of model-free RL \Rightarrow only viable in first order. envs.
- Human being \rightarrow internal model of env dynamics \Rightarrow E. Dupoux: Babies learn mainly from observation, little intervention! \rightarrow reactive time exps.



- 'Filling the blanks' \Rightarrow what the brain does all the time
 \rightarrow prediction as essence of intelligence



- \rightarrow Solutions: Latent Var. Energy-based Models!
 Lo GANs \rightarrow need distance fct. to define meaningful
 Lo Regularized Latent Var EBM \rightarrow optimize latent w.r.t. Euclidean dist.
 Lo but need to restrict latent var!

high energy
 \approx
 low uncertainty

- Video prediction \rightarrow blocks \Leftrightarrow uncertainty
 Lo combine with semantic segmentation
 Lo Mask-RCNN \rightarrow many masks

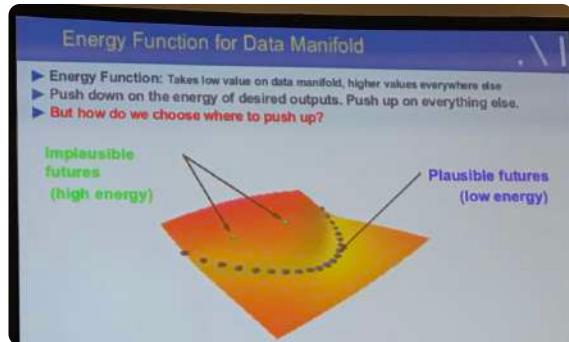
\rightarrow adjoint state method = dual prep
 (control theory) through time

\rightarrow forward model training
 Lo PhysNet ICML 2016
 Lo blurry prediction when overfitting
 \Rightarrow need to uncertainty-weight

average!

average!

average!



- Stochastic latent forward models
 - don't optimize latent, but predict to latent via cerebellum veridical error
 - train policy networks by prop. temporal cost functions \Rightarrow adds uncertainty cost in objective function (MPUR)
 - \hookrightarrow fail \Rightarrow dropout idea!
- no interaction with real world needed \Rightarrow only observation

(2a) One-shot high-fidelity imitation (Papetzel)

- Over-imitation \Rightarrow behavior overfitting \rightarrow kills but not changes
- psychologist belief that human behavior is used to learn new tasks quickly \Rightarrow seems different but is not!
- Def. of imitator research
 - use larger networks; before smaller ones because belief that variance in gradients too high!

(2b) Robotic Manipulation through Visual Planning and Acting (Berkeley)

- Raytracing \Rightarrow visualization of action consequences
- InfoGAN \Rightarrow noise sampled from meaningful latent representation
 - \hookrightarrow causal latent representation \rightarrow second latent vector
- Vision plan \rightarrow inverse model for action taking
- Added many obstacles \rightarrow context-conditioned CIGAN

XMTZLGB
WORK!

(2c) Contingency-based Exploration in RL (Midgley)

- State representability learning for exploration \Rightarrow not only visual similarity
 - \rightarrow Contingency awareness: recognizes that some aspects of env dynamics are under agent's control! \Rightarrow self-location in latent space + awareness
- Attentional dynamics model with spatial attention \rightarrow can target regions of uncertainty
- Exploration - center based with entropically learning

③ 3. Andreas - Linguistic scaffolds for policy learning (Berkeley)

- Text instructions of task in referential environment
 - reward speed acts model \Rightarrow goals / task \rightarrow linguistic pragmatics
 \hookrightarrow quantification with rewards
 \Rightarrow language use \approx gameplay \rightarrow get away from supervised NLP paradigm
 \Rightarrow discr.
 - NLP implications:
 - (1) Train bare listener model
 - (2) Train speaker model that acts directly with bare \rightarrow generator
 \hookrightarrow captions / say desc / instructions follow / generate
 - Problem of outfitting game \rightarrow does not generalize to real language
 - \rightarrow self-play / tree search \Rightarrow within rules of language
 - \uparrow RL for language \uparrow

 \downarrow Language for RL \downarrow
 - learning with policy sketches \Rightarrow verbalization of task as class / context!
 - Options \rightarrow Andreas, Klein & Lester (2014) \Rightarrow only top-level separation via sketches \rightarrow modular composition of sketch policies!
 - Superiority of generalization! \rightarrow structured language gives a lot of help
 - Learning language and acting simultaneously
 \hookrightarrow form of multi-task learning!
 - \hookrightarrow interpretability!
- Answers
to here!

④ S. Kakade - Curiosity and Maximum Entropy Exploration

- RL: algorithms > architectures \Rightarrow importance of saddle points
 \rightarrow exploration via intrinsic motivation (Singh et al (2005))
- Policy implies induced distribution over states \Rightarrow capability of agent's control
 \rightarrow might want high entropy distribution
- \rightarrow don't optimize reward but entropic measures

$$d_{\pi}(s) = (1-\gamma) \sum_{t=1}^{\infty} \gamma^t p(s_t=s | \pi)$$

Curves of state space!

$$H(d_{\pi}) = - \mathbb{E}_{s \sim d_{\pi}} \log d_{\pi}(s) \Rightarrow \pi^* \leftarrow \arg \max_{\pi} H(d_{\pi})$$

\hookrightarrow not concave in d_{π} \rightarrow how to optimize?

- Mixture over distr.: induces stochasticity for optimization

$$d_{\pi_{\text{mix}}}(s) = \sum_i \alpha_i d_{\pi_i}(s)$$

\rightarrow provide proofs!

\rightarrow the inner MDP \Rightarrow KL div.

\hookrightarrow Use other reward factors than simply reward!

\rightarrow Prospective: Exploration from an optimization pov

```

Maximum-entropy policy computation

Initialization:
Set  $\pi_{\text{mix}} = (\alpha_0, C_0)$  with  $\alpha_0 = 1$  and  $C = \pi_0$  to be an arbitrary policy.
for  $t = 0, \dots, T-1$  do
    Compute: the induced state distribution  $d_{\pi_{\text{mix}}}$ .
    Define the reward function  $r_t$  as
         $r_t(s) := \frac{\partial H(d_{\pi_{\text{mix}}})}{\partial d_{\pi_{\text{mix}}}(s)} = -(\log(d_{\pi_{\text{mix}}}(s)) + 1)$ . entropic as for of distr.
    Compute:  $\pi_{t+1} = \text{ApproxPlan}(r_t, c_{\text{approx}})$ .
    Update  $\pi_{\text{mix}} \leftarrow (\alpha_{t+1}, C_{t+1})$  using "learning rate"  $\eta$ :
         $\alpha_{t+1} = ((1-\eta)\alpha_t, \eta)$ . policy that explores next cliff-sticky
end for
return  $\pi_{\text{mix}} = (\alpha_{T+1}, C_{T+1})$ .

```

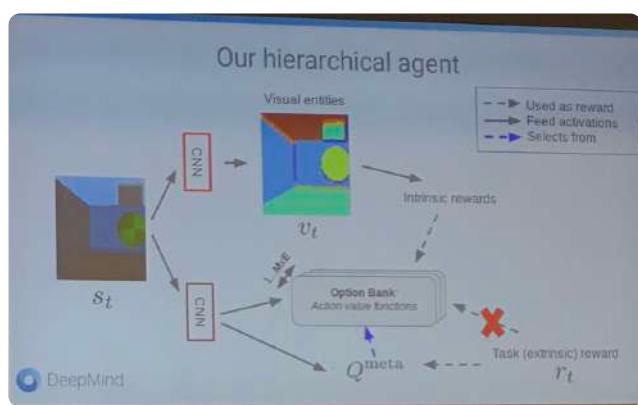
S. M. Kakade (UW)

policy

9/12

Kumar MDP-Course

5a Learning to Control Visual Abstractions for Structured Exploration (DeepMind)



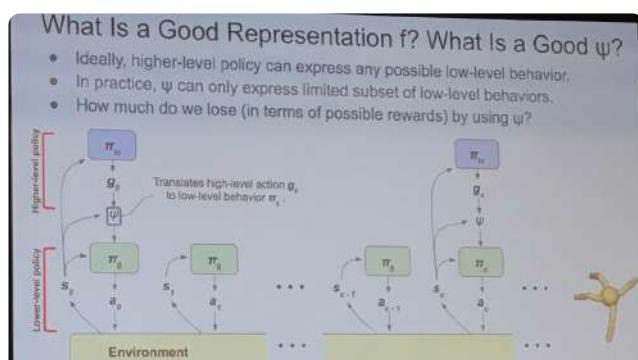
- Visual entities via images extracted by CNN
↳ sequential \Rightarrow contrastive loss
- Extrinsic rewards \Rightarrow geometric
↳ object controllers \rightarrow responsible
↳ not only destroy but different domains (different)

5b Folding Zero : Protein Folding from Scratch (Protein)

- Thermodyn. Hypothesis \rightarrow Hydrophobic - like Model

5c Hier-optimal Representation Learning for HRL (Deep RL)

- Goal-conditioned HRL \rightarrow high-level actions = goals \Rightarrow have to be achieved by low-level
- Problem of defining state representation function ψ
↳ prone over-expansion of constraints



$$\text{SubOpt}(\varphi) = \sup_{S \in S} V^{\varphi}(s) - V^{\varphi_{\text{opt}}}(s)$$

↳ objective fn. with auxiliary fn. approx.
↳ empirical comparison to oracle

⑥ Dimer recap - Knowledge Representation for RL

- Want agent that continually gains its knowledge → no catastrophic forgetting!

$$\mathbb{E} \left[\sum_{t=0}^{\infty} r_{t+1} \prod_{i=t+1}^T \gamma^i | S_t = s, \text{thus} \right]$$

Which representation?

bulky blocks for knowledge

Generalized Value Function!

restrictive: fresh, distant, reward!

relative to UVFT!

relative to successor features
↳ Barreto et al 17, 18

Knowledge representation: Generalized Value Functions (GVFs)

- Given a cumulant function c , state-dependent continuation function γ and policy π , the Generalized Value Function $v_{\pi, \gamma, c}$ is defined as:

$$v_{\pi, \gamma, c}(s) = \mathbb{E} \left[\sum_{k=t}^{\infty} C_{k+1} \prod_{i=t+1}^k \gamma(S_i) | S_t = s, A_{t+\infty} \sim \pi \right]$$

- Cumulant** c is a function that can depend on the S_t , or (S_t, A_t) and can output a vector (even a matrix)
Eg $C_{k+1} = c(S_k, A_k, S_{k+1})$
- Continuation function** γ maps states to $[0, 1]$
- Cf. **Horde architecture** (Sutton et al, 2011); Adam White's thesis; inspiration from Pandemonium architecture

DeepRL workshop 2010

↳ parallel predictor
cephebites

new fresh
↳ can easily
get classic
Value fct.

→ linearity allows us
to combine different VFs!

Rewriting Cumulant +
Continuation function

Successor states and successor features are GVF

- Successor features** (Barreto et al, 2017, 2018) are a natural extension of successor states (Dayan, 1992)
- If states are defined by a feature vector $\phi(s)$, successor features give the discounted sum of future feature vectors from a state.
- In GVF terms, the cumulant is $c = \phi$, and there is a fixed policy and discount
- Interesting property highlighted in Barreto et al:

$$v_{w^T c, \pi, \gamma}(s) = w^T v_{c, \pi, \gamma}(s)$$

which leads to very efficient use of existing GVF to do one-shot computation of new GVF

Option models

- Option model** has two parts:
 - Expected reward** $r_\omega(s)$: the expected return during the execution of option ω from s
 - Needed because it is used to update the agent's internal representations
 - Transition model** $P_\omega(s'|s)$: a sub-probability distribution over next states (reflecting the discount factor γ and the option duration) given that ω executes from s
 - P specifies **where** the agent will end up after the option/program execution and **when** termination will happen
- Models are **predictions** about the future, conditioned on the option being executed

→ think not only from execution point of view!
↳ model ⇒ prediction about future!

The anatomy of the reward option model

- Option reward model is defined as:

$$r_\omega(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s, \omega_t = \omega]$$

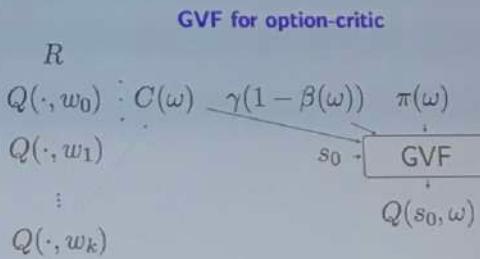
We can write Bellman equations for the model:

$$r_\omega(s) = \sum_a \pi_\omega(s, a) [r(s, a) + \sum_{s'} P_\omega(s'|s) \gamma (1 - \beta_\omega(s')) r_\omega(s')]$$

- This means the **option reward model** is a GVF, $v_{\omega, \gamma}$: cumulant is the environment reward r , policy is π_ω , continuation function is $\gamma(1 - \beta_\omega)$
- Note that we could use $\gamma = 1$ for the overall task and these GVF would still be well defined
- Option transition model can similarly be shown to be a GVF

Option function "slopes"
breaks into pieces

Option as generalized value fct.



- Option-value functions try to do a lot! Which explains some of the pathological behavior observed in option-critic
- Models of options insulate experience within the option
- Intuitively, models should be easier to both learn and use than option-value functions.

↳ stronger / continual learning effect

→ make off-policy more efficient
 → don't update at every time step!

→ Can also rewrite policy gradient pt. 1?
 ↳ mix & match \Rightarrow fine policy
 \Rightarrow GVF allows for easy combination of algorithms \rightarrow L&NOUTOE

↳ Alternatives: Reciprocals

\Rightarrow construct target policy that aligns with behavior!

\Rightarrow variance reduction effect!

• HAL environment choice

\rightarrow structures in env

\rightarrow cont. state-action spaces

These goal solvers are available in discretized version

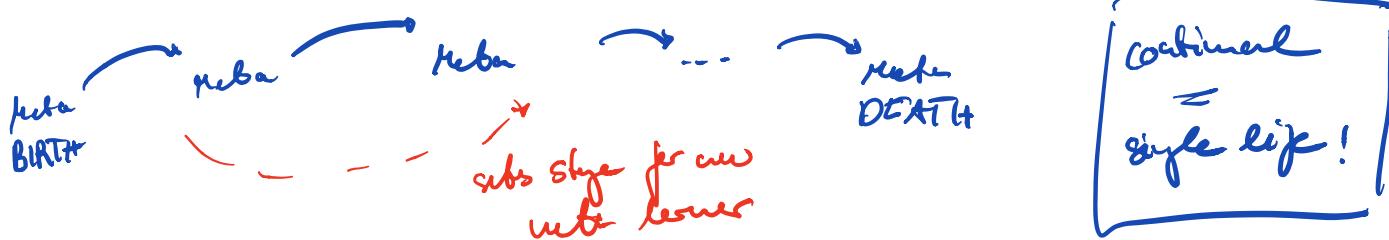
To do list

- How do we generate behavior from GVF?
- Option call-and-return execution gives a hint, but is too limiting
- Ideally, we will use GVF to plan
- Exploration: how to generate data for many GVF?
- Discovery: how do we get initial conditions, policies and continuation functions?
 - Akin to the option discovery problem, but with an emphasis on building models
 - How do we get cumulants?

• Can value func. capture everything relevant to build RL solvers?
 → linear expectation operator \Rightarrow hypothesis needs to be tested!

⊕ 5. Schmidhuber - Circular Lifeacy Meta-heuristics & Artificial Curiosity

- Meta-heuristics + Transfer learning → cost adjustments of few parameters
↳ self-referential, self-modifying \Rightarrow recursive self-improvement
- limits of computability and physics
- critic of Mechanicism \Rightarrow Old 1984 web \rightarrow probabilistic galaxy Net changes itself!

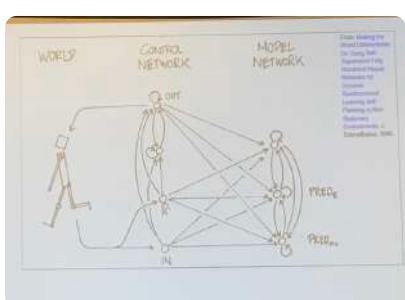


success story AlphaGo \rightarrow set check points \Rightarrow history of record accelerating of self-modification

↳ 'Is life better before or after checkpoint?'

- Gödel Machine (2003) \leftarrow self-modifying

- RNN Controller & World Model \rightarrow controller learns ability to use the model! \rightarrow how to motivate



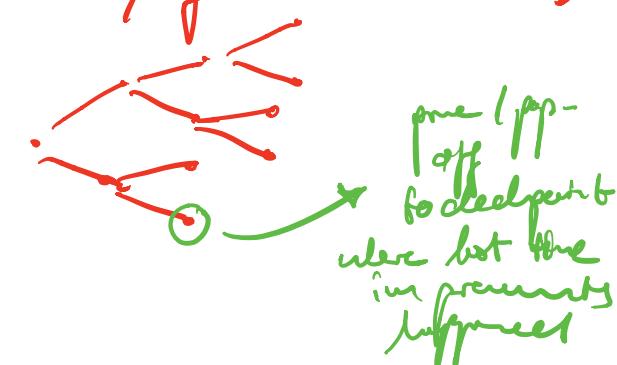
(\rightarrow hyperprior)

Prior

↓

Posterior

- Artificial curiosity \rightarrow adversarial minimax approach
↳ improvement from "play-pool" experiments
↳ informed goals reference



⑧ Y. Gal - Toward's Robust Evaluations of Continual Learning

- continual learning \Rightarrow storing of data \rightarrow adapt to new batches of data
 ↳ E.g. dynamic models, discarding data after learning

Need:

- cross-task resilience
 - general output level
 - no re-training on old tasks
 - few them & tasks
 - Don't know task id at test time \rightarrow which task are we trying to answer at test
 - Real time updates // small memory footprint
- Multi-headed split MIST
 → no catastrophic forgetting
 → 0 vs 1 batch gives add. info!

Why are the used exps bad?

VCL

For each task t :

- Given:
 - Dataset $D_t = \{(x_i^t, y_i^t)\}_{i=1}^{N_t}$
 - posterior from task $t-1$: $\mathcal{N}(\mathbf{w}_{t-1}, \Sigma_{t-1})$
- Optimise:

$$\mathcal{L}_t(\mathbf{w}) := \sum_{x_i, y_i \in D_t} \log p(y_i|x_i, \mathbf{w}) - \frac{\sigma}{2} (\mathbf{w} - \mathbf{w}_{t-1})^T \Sigma_{t-1}^{-1} (\mathbf{w} - \mathbf{w}_{t-1})$$

to obtain posterior for task t : $\mathcal{N}(\mathbf{w}_t, \Sigma_t)$.

¹ EWC, SI, RW similar with Σ_{t-1} def'd differently (e.g. approximate Fisher information).

Posterior Forgets Tasks

Why are the used exps bad?

Coresets

For each task t :

- Given:
 - Dataset $D_t = \{(x_i^t, y_i^t)\}_{i=1}^{N_t}$
 - Coresets C_1, C_{t-1}
- Select coreset $C_t = \{(x_i^t, y_i^t)\}_{i=1}^{N_t}$ (e.g. k-means centroids of D_t).
 Keep for future tasks.
- Optimise:

$$\mathcal{L}_t(\mathbf{w}) := \sum_{x_i, y_i \in D_t} \log p(y_i|x_i, \mathbf{w}) - \frac{\sigma}{2} \mathbf{w}^T \mathbf{w}$$
- When testing for task $t' < t$, retrain with old data:

$$\mathcal{L}_{t'}^{\text{finetune}}(\mathbf{w}; C_{t'}) := \sum_{x_i, y_i \in C_{t'}} \log p(y_i|x_i, \mathbf{w}).$$

Coresets Forgets Tasks

Why are the used exps bad?

VCL with coresets:

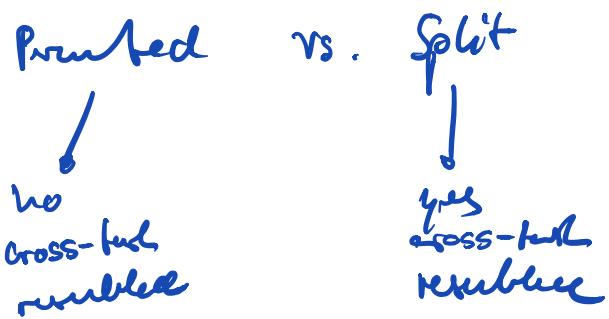
- Optimise

$$\mathcal{L}_t^I(\mathbf{w}) := \sum_{x_i, y_i \in D_t} \log p(y_i|x_i, \mathbf{w}) - \frac{\sigma}{2} (\mathbf{w} - \mathbf{w}_{t-1})^T \Sigma_{t-1}^{-1} (\mathbf{w} - \mathbf{w}_{t-1})$$
- When testing for task t' , retrain with old data

$$\mathcal{L}_{t'}^{\text{finetune}}(\mathbf{w}; C_{t'}) := \sum_{x_i, y_i \in C_{t'}} \log p(y_i|x_i, \mathbf{w}).$$

(Similarly, RW also trains on subsamples of previous datasets, EWC revisits each dataset multiple times)

Combine both ideas!



- Work on single-headed MIST \rightarrow hard experts > hard debuffers
- Discard when task has changed!

⑤ S. Clune - Go-Explore \rightarrow algo for Hard-Exploration Problems

• Hard Exploration problems

\rightarrow sparse-rewards \Rightarrow no

\rightarrow deceptive paths \Rightarrow way \rightarrow local optima

• Intrinsic Motivation \rightarrow Burda et al (2018) \Rightarrow Hartmann's bugs



\rightarrow Detainment

\rightarrow not solved by buffer

• Detainment

\rightarrow first solve, then explore!

\rightarrow robustness effects?

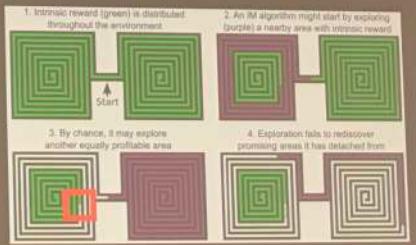
\rightarrow first solve task, then robustify

①

②

Detachment

- Proposal: explicitly remember
 - where promising locations are
 - how to get back to them



\rightarrow all repeatable \Rightarrow similarity after clean-copy

\hookrightarrow phase 2: robustify \Rightarrow add noise!
via backtracks imitation learning

\Rightarrow DISCUSSION STOCHASTICITY TRAIN/TEST

\hookrightarrow 2 modes of behavior!

Backwards Imitation Learning

Starting with demonstration (solution trajectory)

1. Start at the end
2. Backup k steps
3. Run RL until \geq original score
4. go to 1



⑩ bayesian DL funnel

marked before just freefall