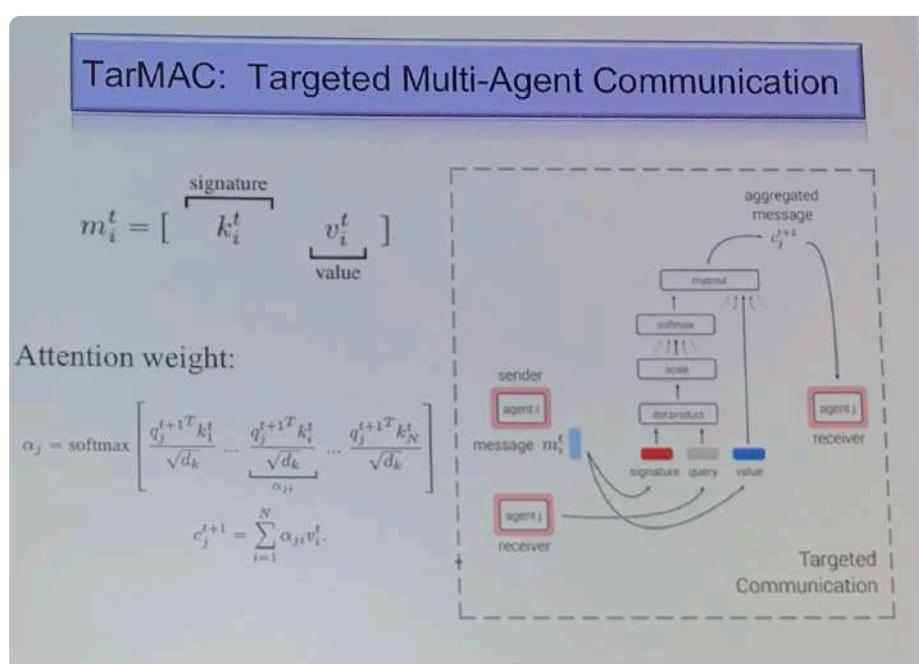
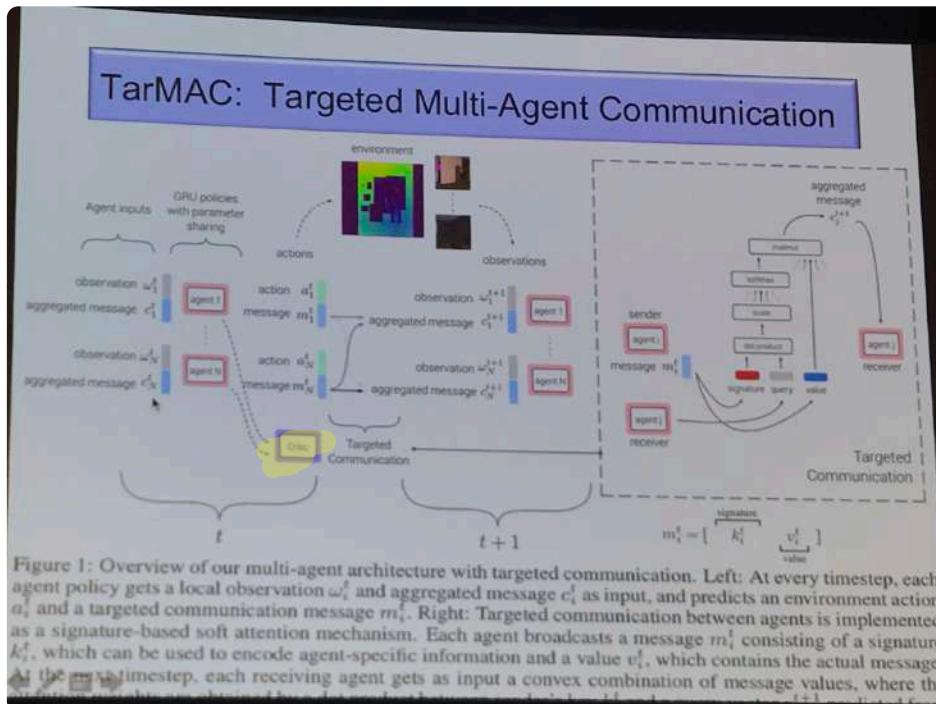


# NIPS - Day 6 - Workshop on SL in Social Obs. Setting

## ① J. Pineau - TarMAC: Targeted Multi-Agent Communication

- cooperative, bandwidth-limited (sending and receiving)!
- $\rightarrow$  no supervisor on communication
- decentralized POMDP setup  $\rightarrow$  agent-specific policy  $\pi_{\theta_i}(\omega_i)$
- with actor-critic setup  $\rightarrow$  centralized critic (joint overall fit.)
- Motivation: self-driving cars  $\rightarrow$  not everyone can communicate at same time



$\hookrightarrow$  if goal is reached by one agent  
 $\hookrightarrow$  spike in communication from one agent to others!

• Trained via backprop through complete architecture

$\rightarrow$  aggregate message: compression of all messages

$\rightarrow$  agent gets localized observation

• message-signature: intend to speak to receiver!

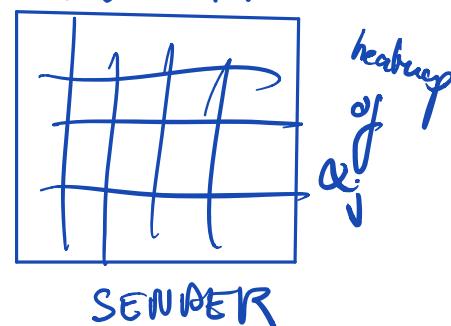
• message-collector: intend to listen to speaker!

• Test in shapes domain

$\hookrightarrow$  agents need own particular goal

$\hookrightarrow$  experiment with causal communication channel  
 $\hookrightarrow$  one agent delays its communication

R  
E  
C  
E  
I  
V  
E  
R



- Also: Traffic junction experiment and have 3D env.
- 1stage / 2stage: multiple rounds of communication  $\rightarrow$  power!
- Critic: architecture not good at planning!  $\Rightarrow$  deep TC method problem

$\rightarrow$  push env. complexity and larger horizons

Comparison to related work				
	Decentralized Execution	Targeted Communication	Multi-Stage Decisions	Reinforcement Learning
DIAL (Foerster et al., 2016)	Yes	No	No	Yes (Q-Learning)
CommNets (Sukhbaatar et al., 2016)	No	No	Yes	Yes (REINFORCE)
VAIN (Hoshen, 2017)	No	Yes	Yes	No (Supervised)
ATOC (Jiang & Lu, 2018)	Yes	No	No	Yes (Actor-Critic)
TarMAC (this paper)	Yes	Yes	Yes	Yes (Actor-Critic)

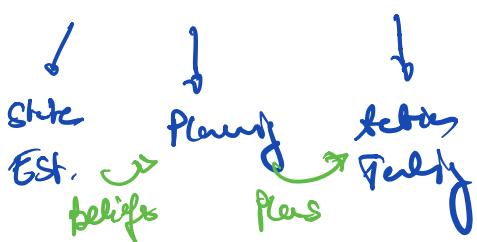
- ② L. Kaelbling - Partial Observability  $\rightarrow$  Physical Robots (MeT)
- Stigmergetic policies (1989)  $\rightarrow$  'Learn' policies with external memory  
 $\rightarrow$  leave traces in env that reward agent of other agents
  - would handle our different envs
    - (1) reverse - eng. humans
    - (2) reciprocate evolution
    - (3) solve POMDP in 'robot policy' of life!
    - (4) optimal policy as RL POMDP algorithm
- general purpose  
↓  
slow
- Biological intelligence  
 $\hookrightarrow$  Spelke & Vinson  
 "Core knowledge"  
 $\hookrightarrow$  invariants in world  
 $\hookrightarrow$  human constraint  
 in agency loop

Build in general representation and inference mechanisms:	
<ul style="list-style-type: none"> <li>• convolution in space and time</li> <li>• kinematics</li> <li>• path planning</li> <li>• forward/backward causal inference</li> </ul>	<ul style="list-style-type: none"> <li>• abstraction over objects</li> <li>• state abstraction/aggregation</li> <li>• temporal abstraction</li> <li>• state estimation, data association</li> </ul>

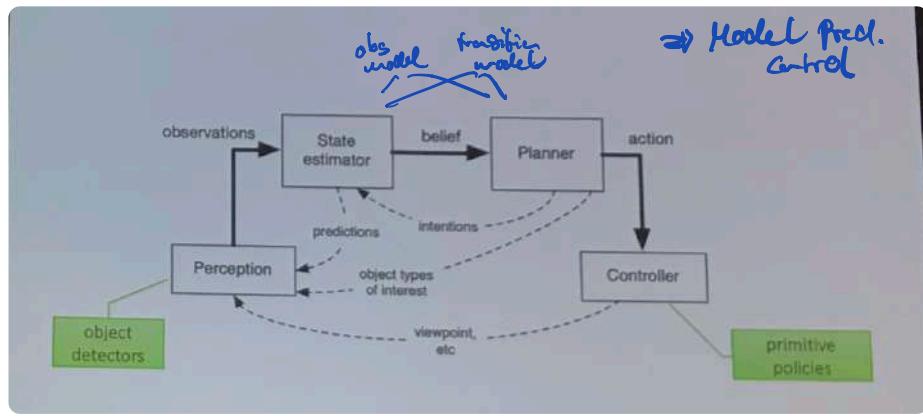
Hard-coded vs. learned

↓  
 learnable models  
 inference rules  
 several actions

WHAT CONSTITUTES RL?

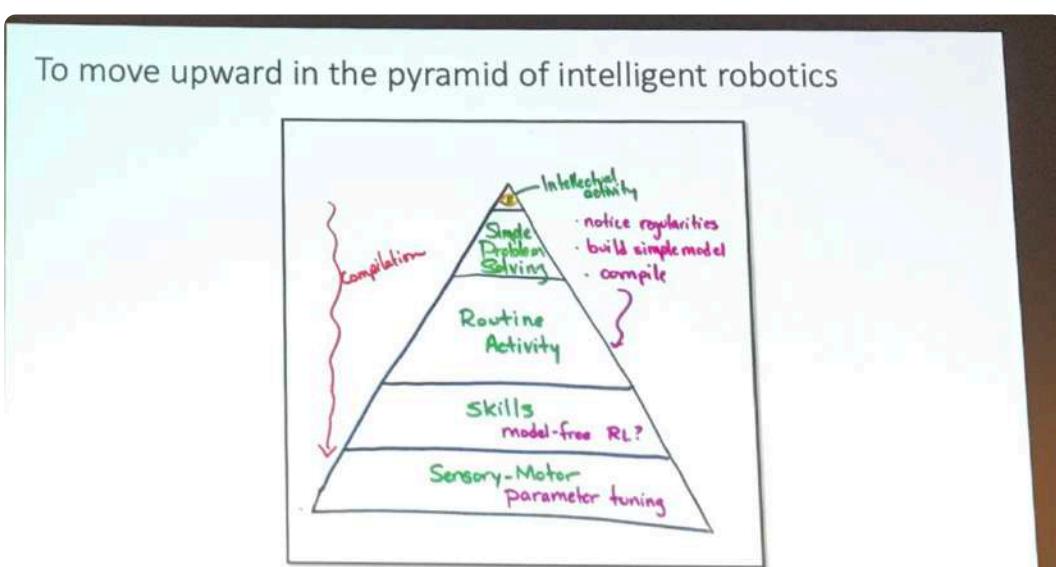


$\Rightarrow$  belief-space model-free control with model adaptation



• Observations ①  
 what action to take → effect on world  
 → effect on my own belief of the world  
 meta-cognition

- Occlusion vs. Noise → epistemic / irreducible noise!  
 ↓  
 hard sys. explanation problem on its own
- Dresden (Made Up Models) books



### (3) I. Gelman - High-level Strategy Sketch (FTIR)

- "Fog of war" = incomplete info
- unknowable task → jointly learn hidden state est. only one effective  
 ↳ e.g., depth estimation → bootstrap learning!  
 ↳ use: predict ground-truth opponent unit acts by type  
 → Optimize parameterless → small MLP on top of LSTM

## (4) D. Silver - Adaptive Deep RL (DeepMind)

- Combat sensitivity of RL to hyperparameters  $\Rightarrow$  can very quickly decay!  
↳ sensitivity to changes in env.  $\Rightarrow$  Online adaptability of meta-params  
↳ within one lifetime!

lose value  
of reserves  
of reward!

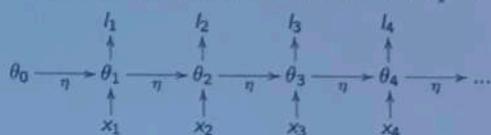
meta = learnable  
hyper = not learnable

### Meta-Parameters

- A deep RL algorithm is defined by update function  $f$  and meta-parameters  $\eta$
  - e.g. TD( $\lambda$ ) adjusts parameters  $\theta$  to minimise prediction error w.r.t. the  $\lambda$ -return
- $$\theta_{k+1} = \theta_k + f(x, \theta_k, \eta)$$
- x: seq. of experiences*
- $$\theta_{t+1} = \theta_t + \alpha \frac{\partial(G_t^\lambda - v_\theta(s_t))^2}{\partial \theta}$$
- The nature of the update is determined by meta-parameters  $\eta = \{\alpha, \gamma, \lambda\}$
  - How can we adapt meta-parameters to optimise an **extrinsic objective**?

### Meta-Gradient Reinforcement Learning (Xu et al.)

Trace effect of meta-parameters on future parameters  $\theta_k$  and extrinsic losses  $l_k$



The gradient of the extrinsic losses w.r.t.  $\eta$  can be accumulated online

$$\theta_{k+1} = \theta_k + f(x, \theta_k, \eta)$$

$$\frac{d\theta_k}{d\eta} \approx \frac{d\theta_{k-1}}{d\eta} + \frac{\partial f(x, \theta_{k-1}, \eta)}{\partial \eta}$$

$$\frac{\partial l_k}{\partial \eta} = \frac{\partial l_k}{\partial \theta_k} \frac{d\theta_k}{d\eta}$$

↳ additive structure  $\rightarrow$  online only  
gradients can also simply be accumulated

- Extrinsic loss = policy gradients

$\rightarrow$  meta-meta gradients?!

$\rightarrow$  need to condition on  
hyperparameters  $\Rightarrow$  as input  
 $\hookrightarrow$  similar UVFT!

- DQN  $\Rightarrow$  overfitted!  $\rightarrow$  help!  
 $\hookrightarrow$  change of problem definition!  
 $\hookrightarrow$  Hessel et al.  $\Rightarrow$  adaptive DQN

### Conditioned Policy and Value Function

- The meta-parameters (e.g.  $\gamma$ ) may determine the **semantics** of the return
- If the meta-parameters change, the semantics of the policy/value change
- It is therefore important to condition the policy and value function on  $\eta$ 
  - Provide  $\eta$  as an input to the policy and value function
  - Learn a "universal" policy/value across different meta-parameters

$$v_\theta^\eta(S) = v_\theta([S; \mathbf{e}_\eta]) \quad \pi_\theta^\eta(S) = \pi_\theta([S; \mathbf{e}_\eta])$$

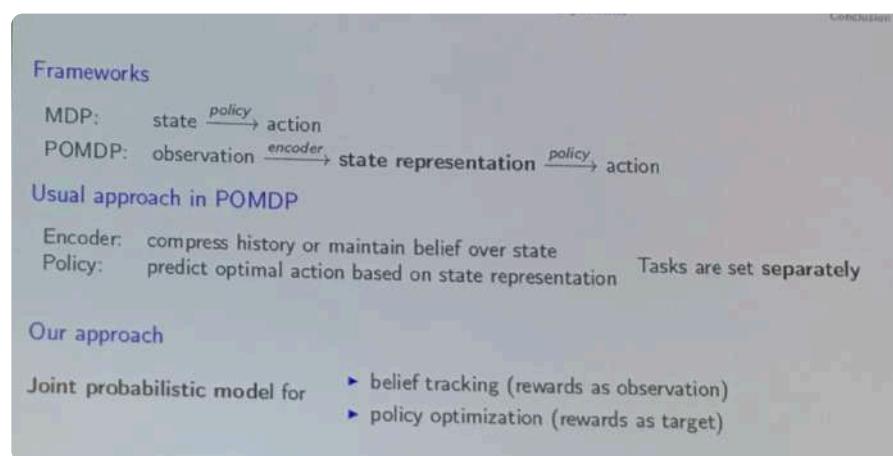
- This was necessary for strong performance

## Learning the RL Algorithm

- Deep RL learns **representations**
  - Value/policy represented by a deep neural network
  - Allowing for flexible and powerful representations
  - Efficiently learned by gradient-based optimisation
- Meta-gradient RL learns **update functions**
  - Update rule can be represented by a deep neural network
  - Allowing for flexible and powerful updates
  - Efficiently learned by gradient-based optimisation
- The ultimate goal:
  - A meta-RL algorithm that learns its own RL update rule
  - Adapted during the agent's lifetime in its environment
  - Meta-gradient cost:  $O(\#parameters \times \#meta-parameters)$

↑ *Utilities scalability!*

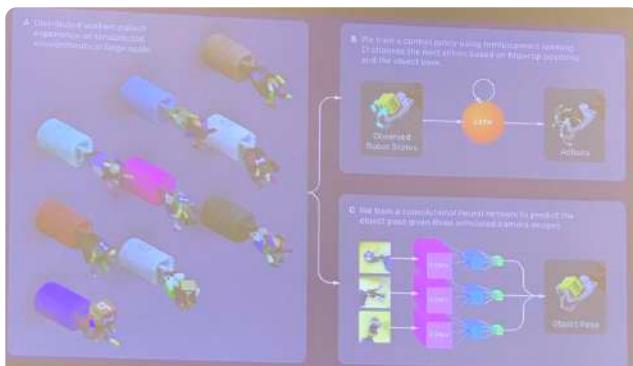
## ⑤ A. Krishin - Joint Belief Policy and Reward Opt. via Approx. Inf.



## ⑥ L. May - Learning Dexterity (OpenAI)

- High-dim control data, hand  $\Rightarrow$  24 dof, 20 actuators  
 $\rightarrow$  partial observability  $\rightarrow$  finger tip views occluded by object!
- Sim2Real: train in simulator and deploy in reality  
 $\rightarrow$  place sparse viewing of finger positions!  
 $\uparrow$   
Physics + Visual! transfer!
- Domain randomization  $\Rightarrow$  trained for robustness!

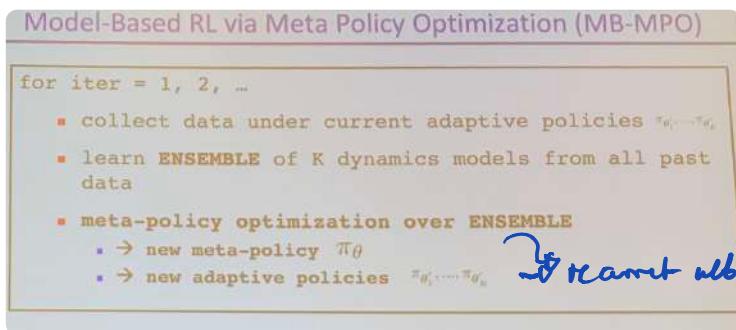
- (1) Dofc Model LSTM
- (2) Curb Model LSTM
- (3) Cision Model  
multi-camera CNN



## ⑦ P. Abbeel - RL under Env. Uncertainty

①

- Model-based  $\Rightarrow$  more sample efficient  $\rightarrow$  but gets to lower level
  - $\hookrightarrow$  overfitting: policy opt. exploits sys. where no support data available  $\Rightarrow$  model-biased  $\rightarrow$  solution: model-ensemble TRPO (MB-TRPO)
- Hindsight: Model-based Meta-Policy Optimization
  - $\hookrightarrow$  adaptive policy robust to dynamics randomization & also robust in real deployment
  - $\hookrightarrow$  learn enough statistics to capture dynamics of the world



- Representation Learning  $\Rightarrow$  for exploration based on prior knowledge ②
  - $\hookrightarrow$  representations of behaviors not necessary static!  $\Rightarrow$  generates!
  - $\hookrightarrow$  Cyriac et al (2018)  $\Rightarrow$  Model-agnostic (AESN)

- Fast RL via few RL
  - $\hookrightarrow$  envs are all mathematical defined  $\Rightarrow$  how to handle advantage of prior dynamics laws
  - $\hookrightarrow$  Merton-hierarchy: learn the learning algo  $\rightarrow$  how do envs need to be like?
  - $\hookrightarrow$  Duan et al (2016): train = test  $\Rightarrow$  every time newer env is generated both rules/ability or problematic!

③

# E. Todorov - Solving estimation and control problems with one unified

Linearly-solvable MDPs

Todorov 2006-2010

In traditional MDPs the controller chooses actions  $u$  which in turn specify the transition probabilities  $p(x'|x, u)$ . We can obtain a linearly-solvable MDP (LMDP) by allowing the controller to specify these probabilities directly:

$x' \sim u(\cdot x)$	controlled dynamics
$x' \sim p(\cdot x)$	passive dynamics
$p(x' x) = 0 \Rightarrow u(x' x) = 0$	feasible control set

The running cost is in the form

$$\ell(x, u(\cdot|x)) = q(x) + D_{\text{KL}}(u(\cdot|x) || p(\cdot|x))$$

Thus the controller can impose any dynamics it wishes, however it pays a price (KL divergence control cost) for pushing the system away from its passive dynamics.

(related results in continuous time)

still can't solve the  
to high-order!  
linear operator

Estimation-control duality

The probability of a trajectory under the passive dynamics is

$$p(x_1, x_2, \dots, x_T | x_0) = \prod_{t=0}^{T-1} p(x_{t+1} | x_t)$$

The probability of a trajectory under the *optimally-controlled* dynamics is

$$\mu(x_1, x_2, \dots, x_T | x_0) \propto \exp\left(-\sum_{k=1}^T q(x_k)\right) p(x_1, x_2, \dots, x_T | x_0).$$

This is equivalent to **Bayesian inference** where  $p$  is the prior,  $\mu$  the posterior, and  $q$  the negative log-likelihood (of some unspecified measurements).

Let  $\mu_k(x)$  be the marginal of the trajectory density at time step  $k$ , and define  $r_k(x) = \mu_k(x) / z_k(x)$ . Then

$$\begin{aligned} z_k(x) &= \exp(-q(x)) \sum_{x'} p(x'|x) z_{k+1}(x') \\ r_{k+1}(x') &= \sum_x \exp(-q(x)) p(x'|x) r_k(x) \\ \mu_k(x) &= r_k(x) z_k(x) \end{aligned}$$

This is equivalent to **recursive Bayesian Inference** where  $r_k$  is the filtering density,  $z_k$  the backward filtering density, and  $\mu_k$  the posterior.

Markov

General case:  
 $x_{t+1} \sim u(\cdot|x_t)$   
 $\ell(x, u) = q(x) + D_{\text{KL}}(u(\cdot|x) || p(\cdot|x))$   
 is dual to  
 $x_{t+1} \sim p(\cdot|x_t)$   
 $y_t \sim \text{Bernoulli}(r(x_t))$   
 when  
 $q(x) = -\log(r(x))$

Linear case:  
 LQR is dual to information filter

Stochastic  
↓  
Det.

↳ solvable since  
for linear systems

- My Ido → based on bad prev. solutions
- bayesian posterior vs. optimal values

Third time is the charm: MDPs (2006-2010)

configuration  
velocity  
applied force  
internal force  
 $M(\cdot)$   
inertia matrix  
 $J(\cdot)$   
constraint Jacobian  
 $\lambda(\cdot, \tau, \tau')$

forward dynamics: inverse optimization

$$\dot{v} = \arg \min_u \| (v + M^{-1}(c - \tau)) \|_M^2 + \gamma (J u - \tau)$$

inverse dynamics: analytical solution

$$\ddot{v} = M \ddot{v} + v + J^T \nabla s(J \ddot{v} - \tau), \quad \lambda = -\nabla s$$

backward dynamics: numerical solution

numerical: 10,100  
numerical: 40,400  
particle: 7,150  
grid: 16,250  
softmax: 33,000

10-core processor

POLO: Plan Online, Learn Online

Indefinite-horizon average-cost Bellman equation:

$$v \rightarrow v(x) = \min_u \{\ell(x, u) + v(f(x, u))\}$$

Equivalent constrained optimization problem:

$$\begin{aligned} \min_u \{v\} \quad \text{s.t.} \quad & (\text{Bellman}) \\ \min_{u \in U} \left\{ \min_u \{\ell(x, u) + v(f(x, u))\} - v(x) \right\} \quad & \text{s.t.} \quad (\text{Bellman}) \end{aligned}$$

Replace the value function  $v(x)$  with an approximation  $\hat{v}(x, \theta)$ .  
 Update the dynamics for  $T$  steps:  $x_{t+1} = f(x_t, u_t)$ .  
 Replace the hard constraint with a soft penalty:

$$\min_{u \in U} \left\{ \sum_{t=0}^T \ell(x_t, u_t) + v(x_T, \theta) - v(x_0, \theta) - c \right\}^2$$

MPG:  $\min_{\{\theta\}} \sum_{\{x\}} \ell(x, u_t) + v(x_T, \theta)$       Learning:  $\min_{\{\theta\}} \sum_{\{x\}} (\ell(x, u_t) + v(x_{T+1}, \theta) - v(x_t, \theta) - c)^2$

lower Rajeswari, Lukac, Todorov and Moritz, arXiv 2010  
so much less samples than policy gradient

- OptiCo SDK
- Costs, derivatives, optimizers

## ③ f. Adams - Inference and Control of Learned Behavior in Robots

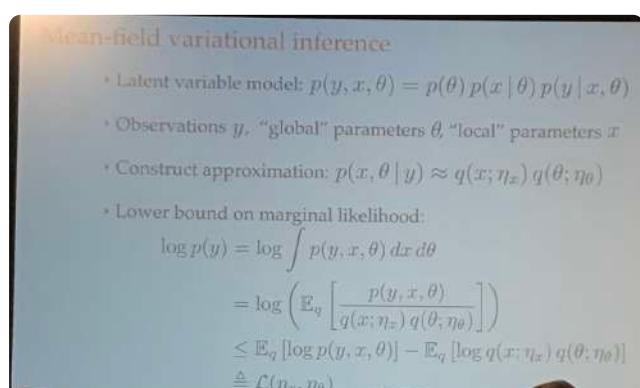
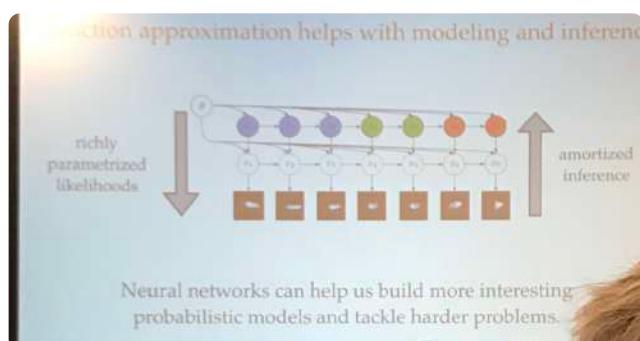
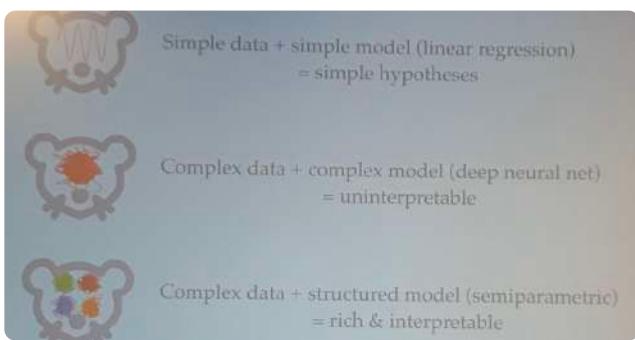
- Task learning in rats (Brody lab micebe)
- high / low tons  $\rightarrow$  2 together with a bag  $\rightarrow$  What is rat's age?
- Roy et al (2018 - NIPS)!  $\Rightarrow$  E-E, disagreement
- observation: slow learning in robots
- Wiltschko et al (2015)



→ unsupervised discovery of language of behavior

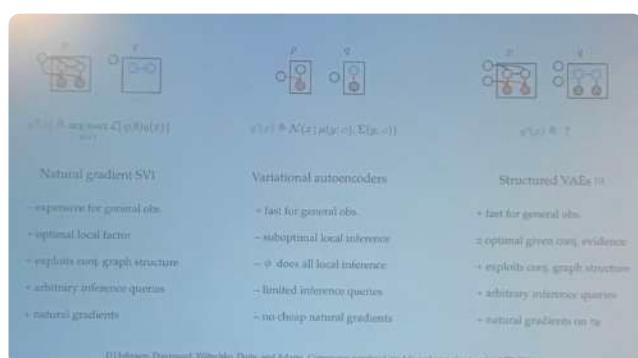
→ mouse = switching LPS  
 $\Rightarrow$  Kalman Filter

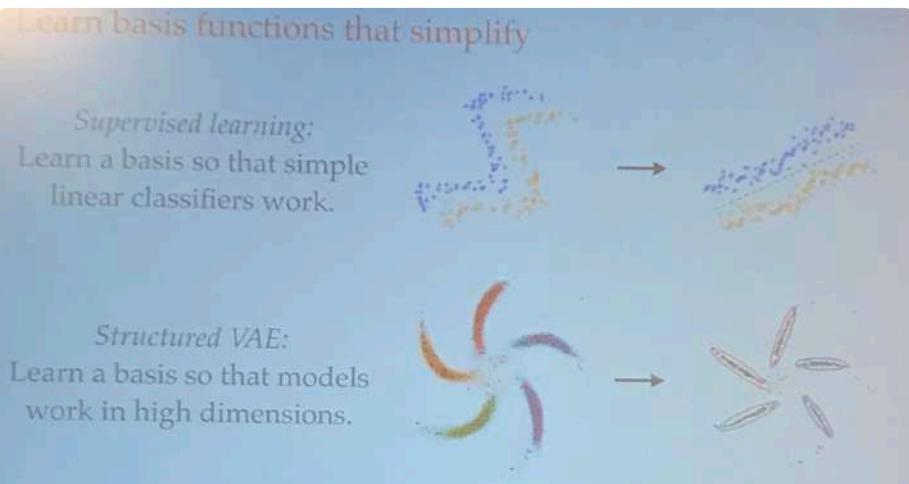
I can also use VAE instead!



→ turn into stochastic VI on global level!

amortized inference for non-conjugate models  
 $\hookrightarrow$  Kingma and Welling  
('Auto-encoding variational Bayes')  
(2014)





→ Fins (2018)  
talk by  
Gordon Bernou

! makes also sense  
for another  
of want  
prior shape!