

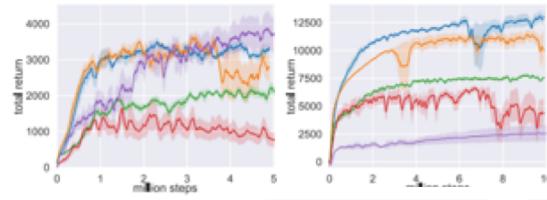
# Hierarchical Reinforcement Learning

Ivy

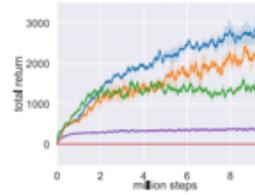
# Outline

- **Motivation and problem Recap (30min)**
- Method (30min)
  - Option-Critic Architecture
  - Feudal Network
- Method in more detail: Latent Space Policies (1h)
  - Maximum Entropy RL
  - Hierarchical training
  - Invertible transformation between levels

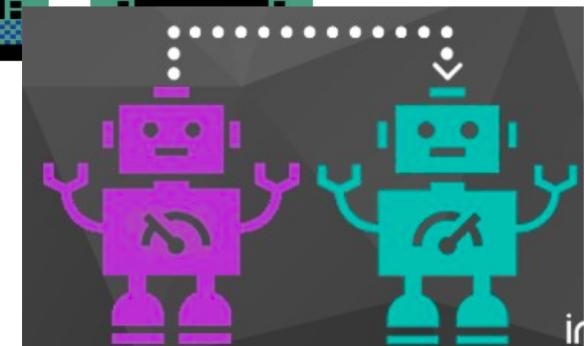
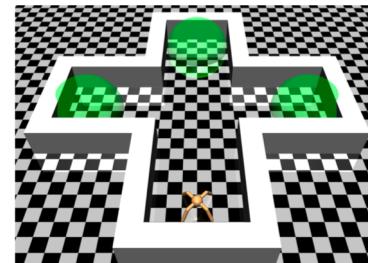
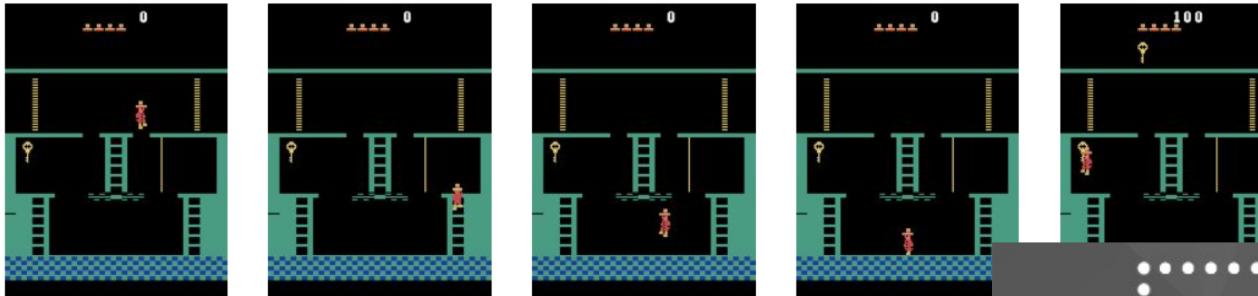
# Recap: Motivating problems



(c) Walker2d-v1



(e) Ant (rllab)



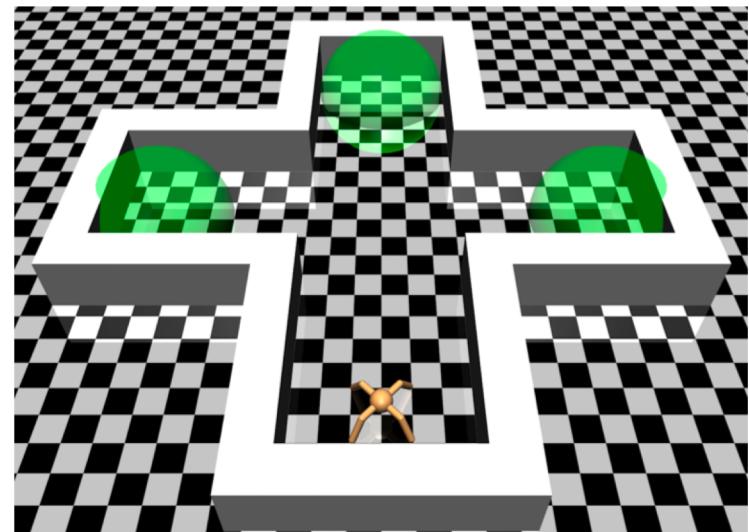
in

# Recap: Sparse reward problem

## Challenges:

- Sparse reward signals
- Long term credit assignment

Temporal abstraction encourages long term credit assignment

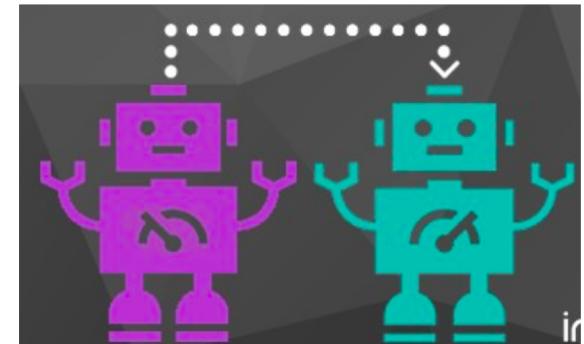


# Recap: Transfer learning problem

## Challenges:

- Disentangle representations

Different levels of the hierarchies can be disentangled.



## Recap: Hierarchy

Hierarchical representation: enable reasoning and decision making at different levels of abstraction

# Recap: Variations of hierarchy

Hierarchy is defined differently:

Could be algorithm being trained, i.e., hierarchical DQN

Could be architecture that has high-level and low-level, i.e., feudal neural network, option-critic

Could be subtask design, i.e., learning and transfer of modulated locomotor controllers

Could be reward design, i.e., value function decomposition

Could be training procedure, i.e., latent space policies for RL

Could be imitation..., i.e., hierarchical imitation for RL

For this talk, we focus on generalizable architectures and follow the line of how they solve the sparse reward problem.

# Recap: Value function, Policy gradient, Actor-critic

Value function

$$V^\pi(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s; \pi\right]$$

Policy gradient

$$\nabla_\theta J(\theta) \approx \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$$

Actor-critic method

$$\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) - \hat{V}_\phi^\pi(\mathbf{s}_i)$$

$$\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$$

# Outline

- Motivation and problem Recap (30min)
- Method (30min)
  - **Option-Critic Architecture**
  - Feudal Network
- Method in more detail: Latent Space Policies (1h)
  - Maximum Entropy RL
  - Hierarchical training
  - Invertible transformation between levels

# Option: Intuition

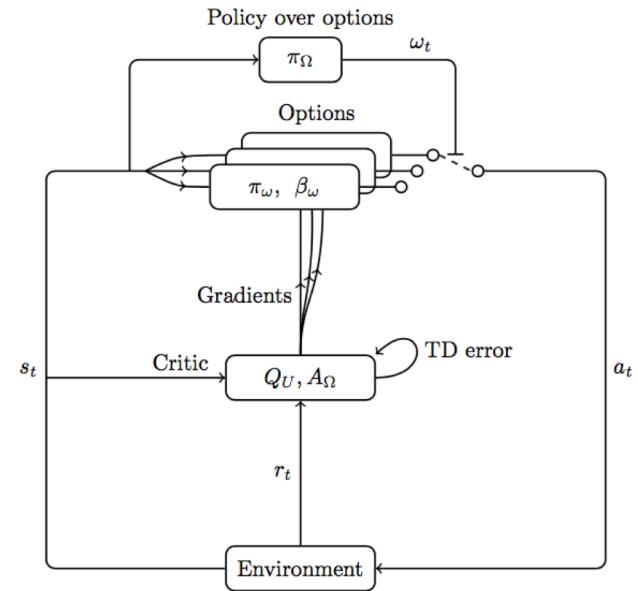
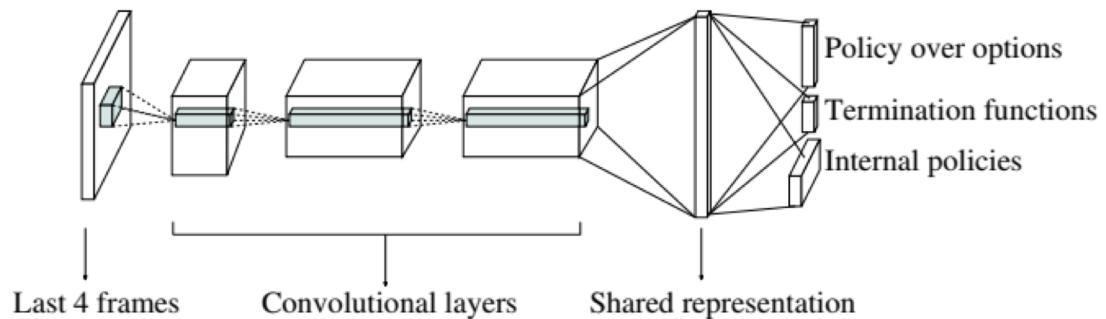
Motivation: Option extended action

Definition: Option is defined as a tuple  $w = \langle \mathcal{I}, \pi, \beta \rangle$

- A policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$
- A termination condition  $\beta : \mathcal{S}^+ \rightarrow [0, 1]$
- An initiation set  $\mathcal{I} \subset \mathcal{S}$

Example: Open-the-door

# Option-Critic: Architecture



# Option-Critic: update

Policy gradient:

$$\nabla_{\theta} J(\theta) = \sum_s d_{\pi_{\theta}}(s) \sum_a \frac{\partial \pi_{\theta}(a|s)}{\partial \theta} Q_{\pi_{\theta}}(s, a)$$
$$d_{\pi_{\theta}}(s) = \sum_t \gamma^t P(s_t = s | s_0)$$

Intra option policy gradient:

$$\nabla_{\theta} J(\theta, w) = \sum_{s,w} d_{\pi_{\theta}}(s, w) \sum_a \frac{\partial \pi_{w,\theta}(a|s)}{\partial \theta} Q_{\pi_{w,\theta}}(s, w, a)$$
$$d_{\pi_{\theta}}(s, w) = \sum_t \gamma^t P(s_t = s, w_t = w | s_0, w_0)$$

# Option-Critic: Experiment

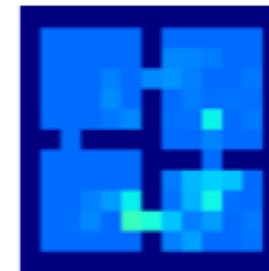
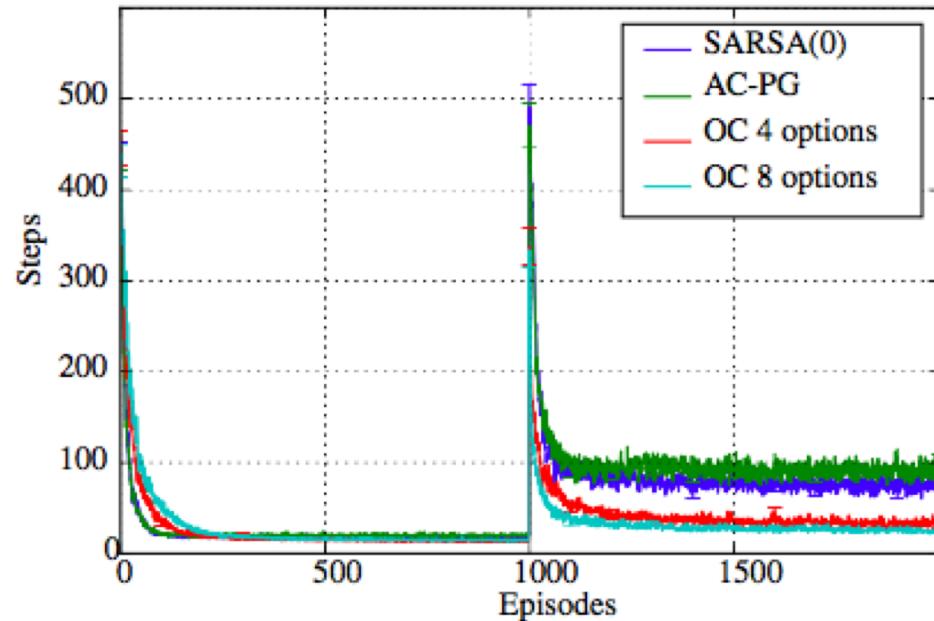


Fig1. Four-room maze task where the target moves after every 1000 episodes. A good algorithm takes the fewest steps to recover the changed target. Option-critic ("OC") recovers faster than actor-critic and SARSA(0).

# Option-Critic: Experiment

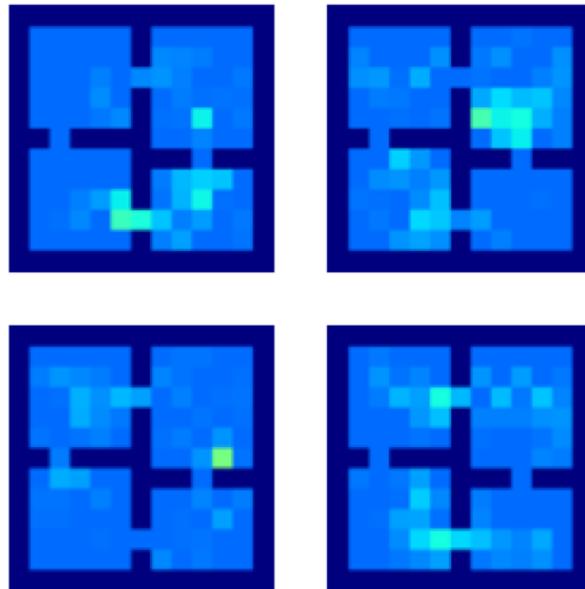


Fig 2. Four-room maze task, Termination probabilities for the option-critic agent learning with 4 options. The black color is the walls of the maze. Lighter colors encode higher termination probabilities.

# Option-Critic Intuition for Solving Sparse Reward problem

- Temporal abstraction: option (extended action).

# Option-Critic: Critique

## Benefits:

- End-to-end
- Learns at speed comparable or better than using just primitive actions

## Drawbacks:

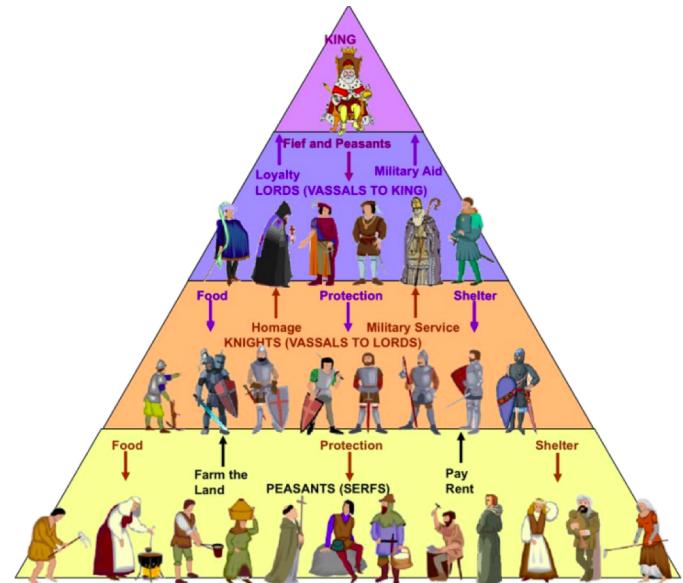
- Tend to converge to trivial solutions as all policy gradient algorithms, need entropy regularization

# Outline

- Motivation and problem Recap (30min)
- Method (30min)
  - Option-Critic Architecture
  - **Feudal Network**
- Method in more detail: Latent Space Policies (1h)
  - Maximum Entropy RL
  - Hierarchical training
  - Invertible transformation between levels

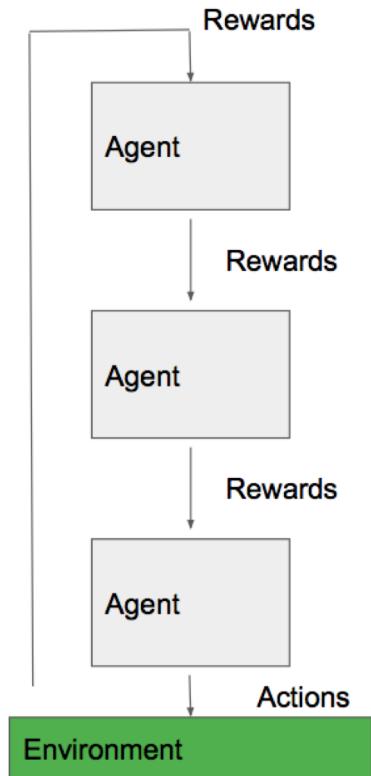
# FuN: Intuition

- Governance system in Europe in middle ages
- Extremely-hierarchical based on ownership of property
- Higher level people have control of the lower levels, but not over people many layers lower.



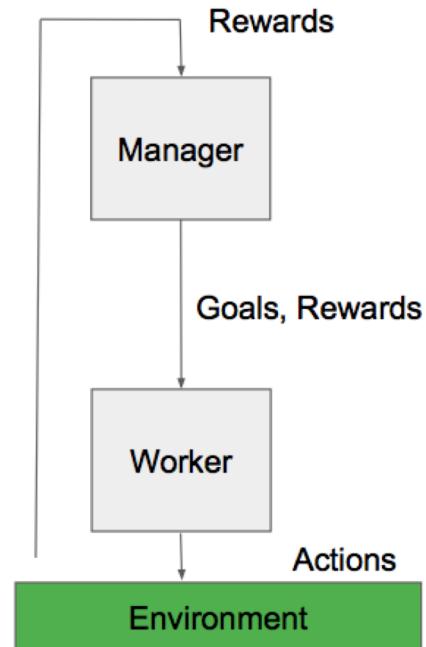
# FuN: Intuition (Feudal RL paper)

- Reward Hiding:
  - Managers reward sub-managers for satisfying their commands, not through an external reward
  - Managers have absolute control
- Information Hiding
  - Observe world at different resolutions
  - Managers don't know what happens at other levels of the hierarchy



# FuN: Intuition

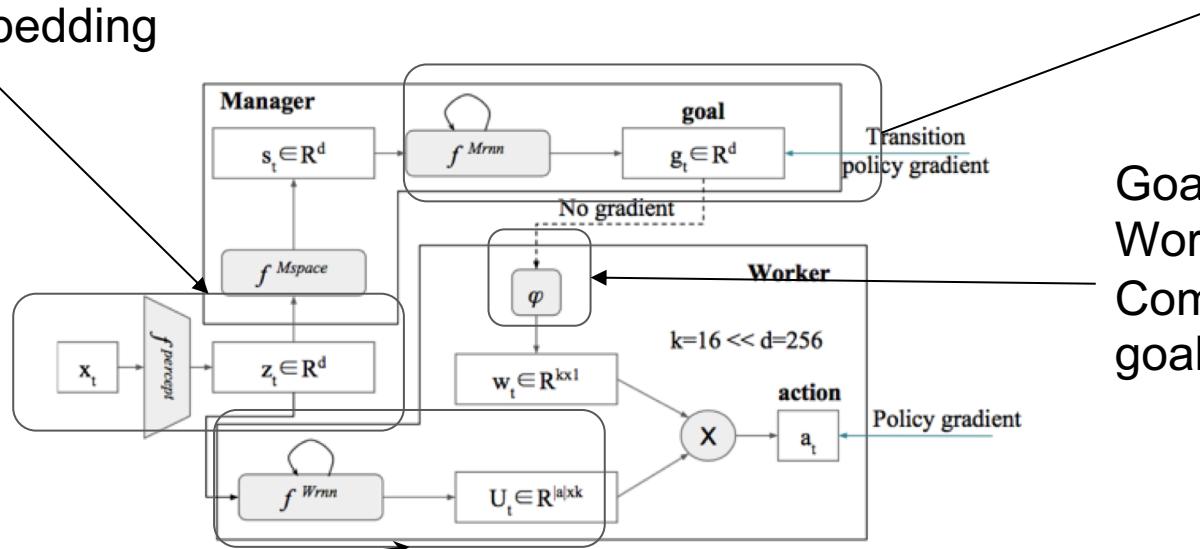
- Manager
  - Sets directional goals for the worker
  - Rewarded by environment
  - Does not directly act in environment
- Worker
  - Higher temporal resolution
  - Reward for achieving manager's goals
  - Produces primitive actions in environment



# FuN: Architecture

Shared visual  
input embedding

Manager's Goal embedding  
(goals summed over last 10 time steps)



Worker's Action embedding

# FuN: Manager's update rule

- Directional goals

$$\nabla g_t = A_t^M \nabla_{\theta} d_{\cos}(s_{t+c} - s_t, g_t(\theta))$$

$$A_t^M = R_t - V_t^M(x_t, \theta)$$

M: manager  
Rt: environment reward

# FuN: Worker's update rule

- Update policy

$$\nabla \pi_t = A_t^D \nabla_{\theta} \log \pi(a_t | x_t; \theta)$$

$$A_t^D = (R_t + \alpha R_t^I - V_t^D(x_t; \theta))$$

D: worker

R<sub>tl</sub>: intrinsic reward

# FuN: Experiments

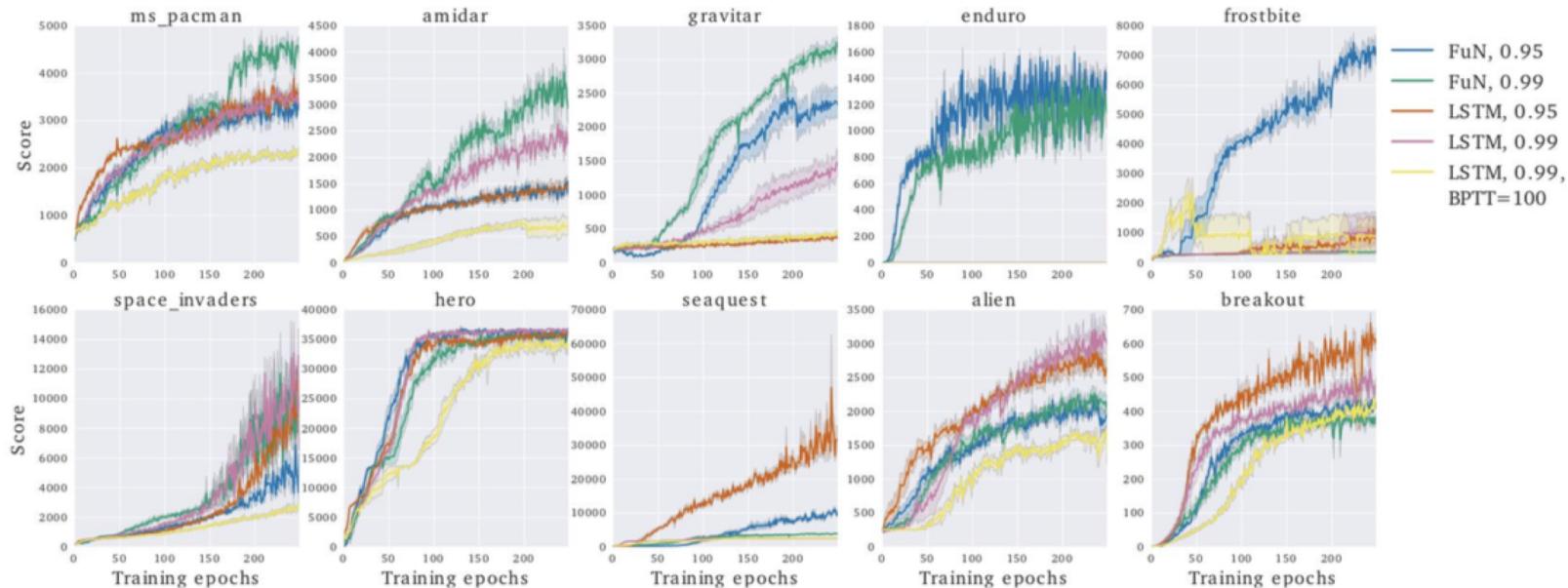


Fig 3, Atari games, compared with baselines.

# FuN: Experiments

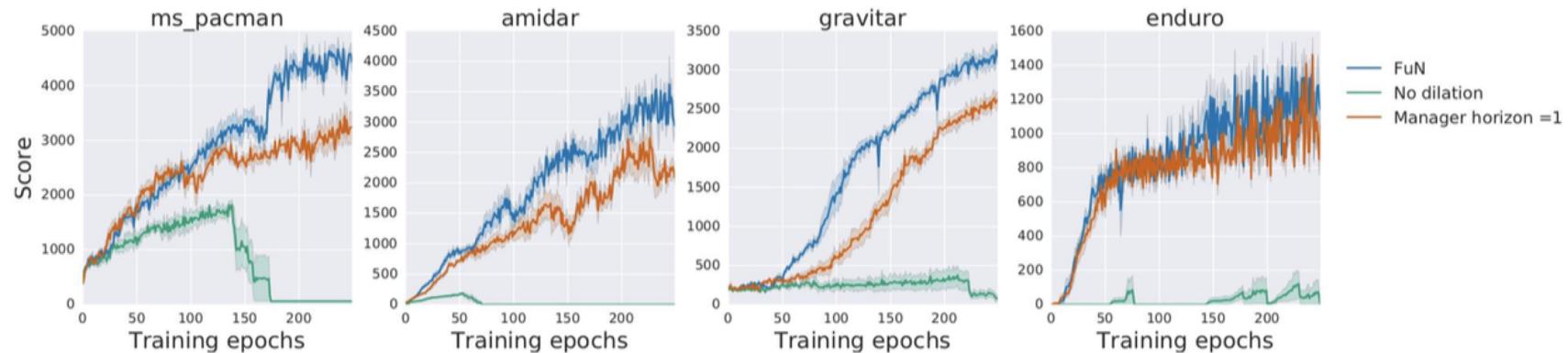


Fig 3, Ablation on temporal resolution ablations

# FuN intuition for solving the sparse reward problem

- Temporal abstraction: enforce lower temporal resolution through dilated LSTM and summing goal over 10 timesteps
- Decouple goal setting and goal achievement: Manually separated manager and worker to make sure they are not learning the same thing.
- Subgoal setting: directional goals rather than absolute goals

# FuN: Comparison with Option-Critic

## Option-Critic

- Adding entropy regularization to avoid trivial policies
- Process-over-options and options are trained together
- Based on option (extended action)

## FuN

- Subgoals as directional goals are naturally diverse.
- Manager and Worker are trained separately with different goals
- Based on intrinsic reward from the manager

# FuN: Critique

Directional goal is a subgoal design

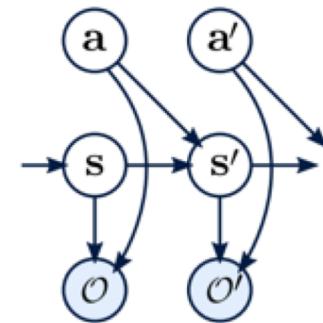
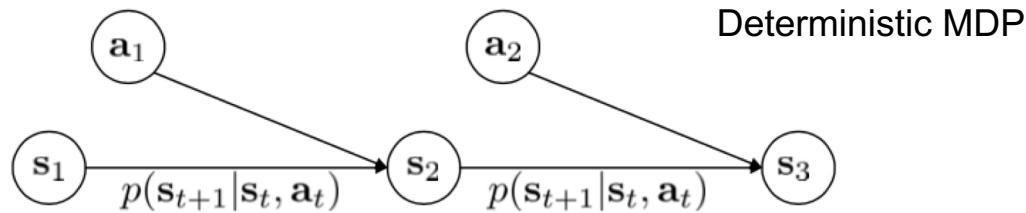
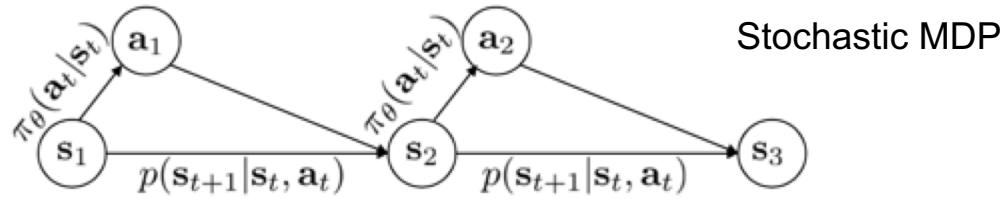
In problems such as locomotion, good subgoal design is hard to come by.

For locomotion task, can we learn good subgoals?

# Outline

- Motivation and problem Recap (30min)
- Method (30min)
  - Option-Critic Architecture
  - Feudal Network
- **Method in more detail: Latent Space Policies (1h)**
  - Maximum Entropy RL
  - Hierarchical training
  - Invertible transformation between levels

# Maximum Entropy RL and optimal trajectory



Set  $p(\mathcal{O}_t|s_t, \mathbf{a}_t) = \exp(r(s_t, \mathbf{a}_t))$

# Maximum Entropy RL and optimal trajectory

$$p(\tau) = p(\mathbf{s}_1, \mathbf{a}_t, \dots, \mathbf{s}_T, \mathbf{a}_T | \theta) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{a}_t | \mathbf{s}_t, \theta) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t).$$
$$p(\mathcal{O}_t | \mathbf{s}_t, \mathbf{a}_t) = \exp(r(\mathbf{s}_t, \mathbf{a}_t))$$

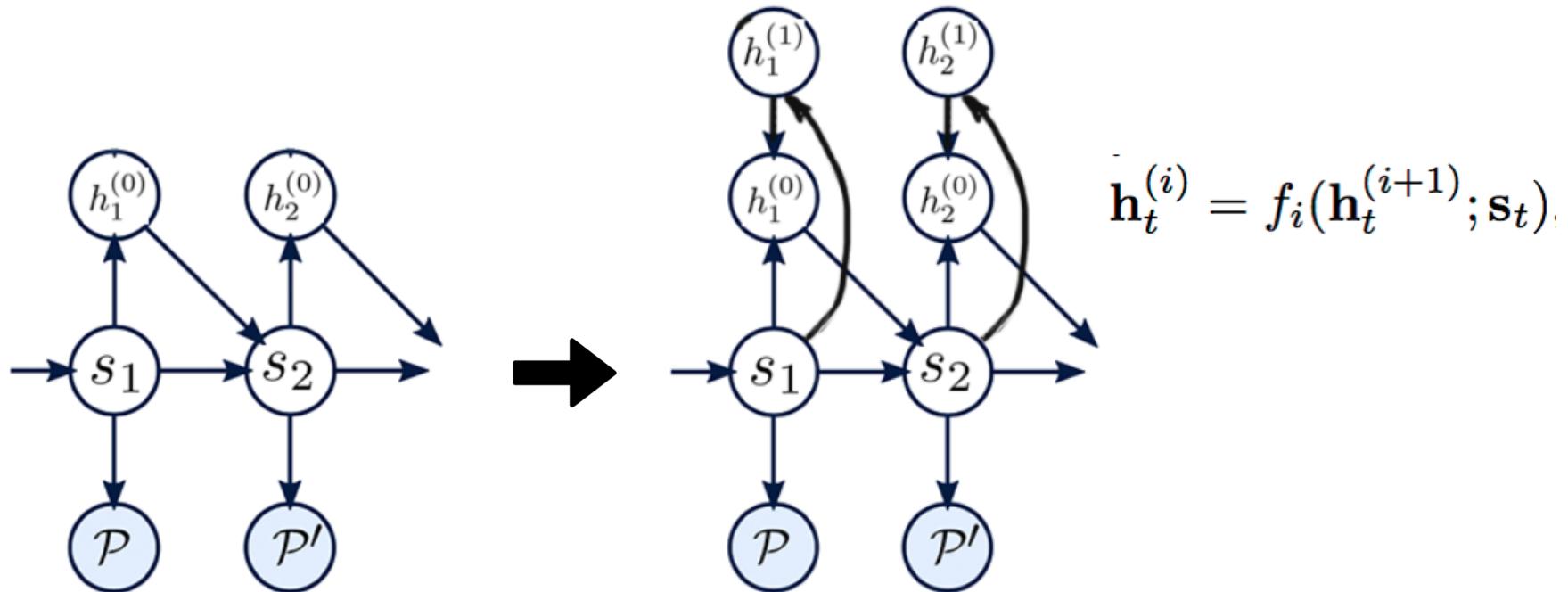
$$p(\tau | \mathcal{O}_{0:T}) \propto p(\mathbf{s}_0) \prod_{t=0}^T p(\mathbf{a}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \exp(r(\mathbf{s}_t, \mathbf{a}_t))$$

Thus, to approximate the optimal trajectory with  $\hat{p}(\tau)$ , minimize  $D_{KL}(\hat{p}(\tau) || p(\tau | \mathcal{O}_{0:T}))$

And since  $-D_{KL}(\hat{p}(\tau) || p(\tau)) = \sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim \hat{p}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t) + \mathcal{H}(\pi(\mathbf{a}_t | \mathbf{s}_t))]$

⇒ The optimal trajectory maximizes the maximum entropy objective

# LSP: Training - Intuition



# LSP: Training a Multi layer policy

**Input:** True environment  $p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{h}_t^{(0)})$ , where  $\mathbf{h}_t^{(0)}$  corresponds to the physical actions  $\mathbf{a}_t$ .

**Input:** Reward  $\mathcal{R}$

**for**  $i = 0$  to  $K - 1$  **do**

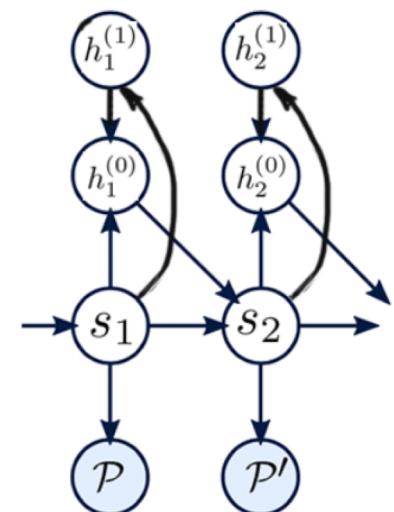
Learn the weights of  $f_i$  so that  $\mathbf{h}_t^{(i)} = f_i(\mathbf{h}_t^{(i+1)}; \mathbf{s}_t)$ ,  
where  $\mathbf{h}_t^{(i+1)} \sim p(\mathbf{h}_t^{(i+1)})$ , optimizes  $\mathcal{R}$   
on  $p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{h}_t^{(i)})$ .

Embed the new layer  $f_i$  into the environment:

$$p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{h}_t^{(i+1)}) \leftarrow p(\mathbf{s}_{t+1} | \mathbf{s}_t, f_i(\mathbf{h}_t^{(i+1)}; \mathbf{s}_t)).$$

**end for**

**Output:** A hierarchical policy  $f = f_0 \circ f_1 \circ \dots \circ f_{K-1}$ .

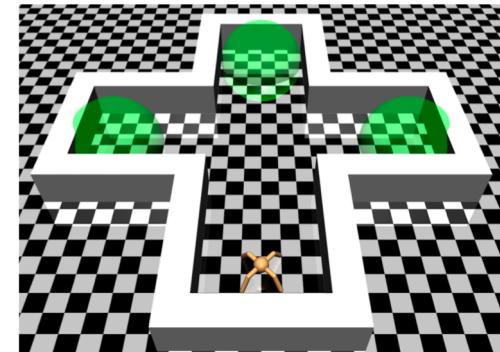


# LSP: Training a Multi layer policy w/ Reward Design

Example (spider navigation):

R1: Low-level reward as motion in any direction, still aim to maximum entropy  $\Rightarrow$  motion in many different directions

R2: final reward, i.e., navigation to target



# LSP: Invertible transformation between levels - Intuition

$$\mathbf{h}_t^{(i)} = f_i(\mathbf{h}_t^{(i+1)}; \mathbf{s}_t)$$

- The higher layer can undo any behavior in the lower layer if it is detrimental to task success.
- Both higher layer and lower layer retain full expressivity

# LSP: Invertible transformation between levels

- Representation of sub-policies:
  - Each subpolicy should be tractable
  - The subpolicies should be expressive, i.e., not cripple the ability of the higher level to solve the task.
  - Conditional factor of subpolicies should be deterministic.

⇐ Real NVP neural network

# LSP: Experiments

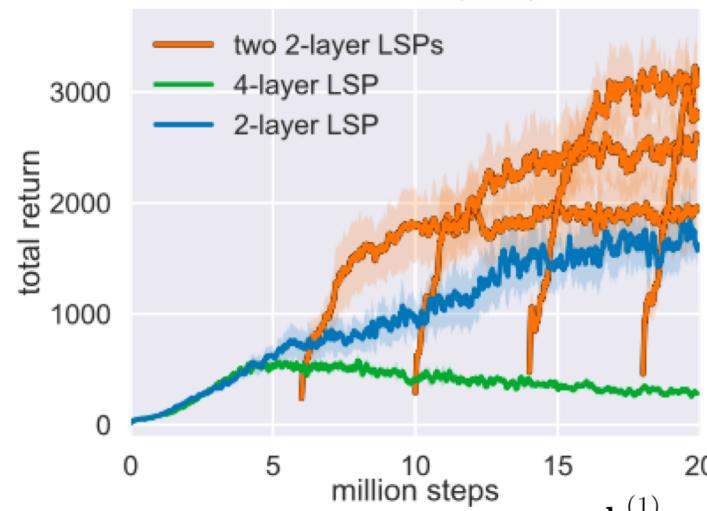
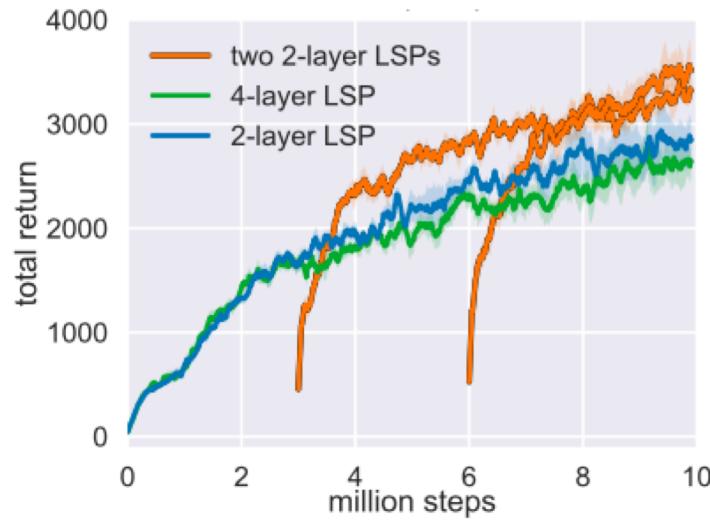


Fig 3. Ant and Humanoid tasks, Single policy with 2 layers, single policy with 4 layers, and two-level policy. Training of two-level policies yields the best performance.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

# LSP: Experiments

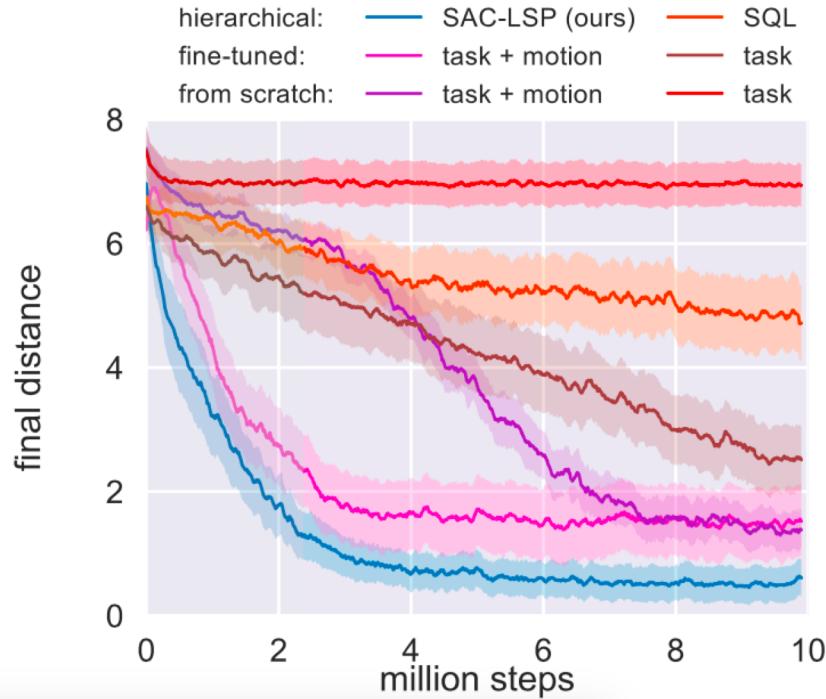


Fig 4, Spider navigation task. We first trained a low-level policy with a motion reward in a pretraining environment not including the walls, then fixed the policy and trained another policy level with a target reward

# LSP Intuition for Solving Sparse Reward problem

- Each layer directly attempts to solve the maximum entropy objective and thus help the layer above it to solve the same maximum entropy objective.
- Bijective transformation - high level and low level are not constrained, but the high level can always reverse what low level did.

# LSP: Comparison with FuN

Feudal NN:

- Diverse set of goals by learning directional goals.
- High level action disintegrated from low level action.
- Trained with different rewards

Option-Critic:

- Entropy regularization
- Extended action.
- Trained together

LSP:

- Maximum entropy objective
- Latent action space.
- Trained level by level

Thank you!

## Control as probabilistic inference: deterministic

$$p(\tau) = p(\mathbf{s}_1, \mathbf{a}_t, \dots, \mathbf{s}_T, \mathbf{a}_T | \theta) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{a}_t | \mathbf{s}_t, \theta) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t).$$

$$p(\tau | \mathbf{o}_{1:T}) \propto p(\tau, \mathbf{o}_{1:T}) = \left[ p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right] \exp \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right)$$

$$p(\tau | \mathbf{o}_{1:T}) \propto \mathbb{1}[p(\tau) \neq 0] \exp \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right)$$

# Control as probabilistic inference: deterministic

In case of exact inference

$$D_{\text{KL}}(\hat{p}(\tau) \parallel p(\tau)) = 0$$

$$D_{\text{KL}}(\hat{p}(\tau) \parallel p(\tau)) = -E_{\tau \sim \hat{p}(\tau)}[\log p(\tau) - \log \hat{p}(\tau)]$$

# Control as probabilistic inference: deterministic

$$\begin{aligned} -D_{\text{KL}}(\hat{p}(\tau) \| p(\tau)) &= E_{\tau \sim \hat{p}(\tau)} \left[ \log p(\mathbf{s}_1) + \sum_{t=1}^T (\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + r(\mathbf{s}_t, \mathbf{a}_t)) - \right. \\ &\quad \left. \log p(\mathbf{s}_1) - \sum_{t=1}^T (\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \log \pi(\mathbf{a}_t | \mathbf{s}_t)) \right] \\ &= E_{\tau \sim \hat{p}(\tau)} \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t) \right] \\ &= \sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim \hat{p}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)] \\ &= \sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim \hat{p}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)] + E_{\mathbf{s}_t \sim \hat{p}(\mathbf{s}_t)} [\mathcal{H}(\pi(\mathbf{a}_t | \mathbf{s}_t))]. \end{aligned}$$

$$-D_{\text{KL}}(\hat{p}(\tau) \| p(\tau)) = \sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim \hat{p}(\mathbf{s}_t, \mathbf{a}_t))} [r(\mathbf{s}_t, \mathbf{a}_t) + \mathcal{H}(\pi(\mathbf{a}_t | \mathbf{s}_t))]$$

# Control as probabilistic inference: stochastic

$$p(\tau) = \left[ p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right] \exp \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right)$$

$$\begin{aligned} \log p(\mathcal{O}_{1:T}) &= \log \int \int p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\ &= \log \int \int p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \frac{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\ &= \log E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[ \frac{p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \right] \\ &\geq E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} [\log p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) - \log q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})] \end{aligned}$$

$$\begin{aligned}
\log p(\mathcal{O}_{1:T}) &= \log \int \int p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\
&= \log \int \int p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \frac{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\
&= \log E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[ \frac{p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \right] \\
&\geq E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} [\log p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) - \log q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})]
\end{aligned}$$

$$\log p(\mathcal{O}_{1:T}) \geq E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) - \log q(\mathbf{a}_t | \mathbf{s}_t) \right]$$