# Interactive Machine Learning: Robotics



### Problems with Robots

- Robots operate in a world they cannot properly model
  - All models are wrong, some models are useful
- The data that a robot's algorithms depend upon are noisy, non-iid, and occasionally non-stationary
- Human inputs help compensate for this
  - Humans are also non-stationary data sources
  - Humans also bias the samples they provide
  - Humans don't share preferences or strategies







## Tutorial Goals for this section

- Gain intuition for how interactive machine learning is used in robotics
- Build familiarity with the terms and techniques used in the field
- Communicate important ideas for making robots useful in practice

- Build a deep understanding of statistical methods at the core of leading learning from demonstration methods
- Provide implementation-level detail for these techniques
- Walkthrough correctness or convergence proofs





### Algorithmic Human-Robot Interaction

#### • Acquiring Skills and Tasks from Demonstration

- Trajectories and Keyframes for Kinesthetic Teaching: A Human-Robot Interaction Perspective
- Learning and Generalization of Complex Tasks from Unstructured Demonstrations
- Autonomously Constructing Hierarchical Task Networks for Planning and Human-Robot Collaboration
- Towards Robot Adaptability in New Situations

#### Cooperative Task Execution

- Interpretable Activity Recognition
- Cooperative Inverse Reinforcement Learning
- Game-Theoretic Modeling of Human Adaptation in Human-Robot Collaboration
- Effective Robot Teammate Behaviors for Supporting Sequential Manipulation Tasks
- Improving Robot Controller Transparency Through Autonomous Policy Explanation

#### • Interaction Design:

Designing Interactions for Robot Active Learners

#### Activity **Recognition** Workflow



#### Activity Generation Workflow



#### Trajectories and Keyframes for Kinesthetic Teaching: A Human-Robot Interaction Perspective

[HRI 2012]

Baris Akgun, Maya Cakmak, Jae Wook Yoo, Andrea L. Thomaz

Trajectories and Keyframes for Kinesthetic Teaching: A Human-Robot Interaction Perspective

- Multiple methods exist for skill learning on a robot
- Kinesthetic teaching removes the correspondence problem
- When is it appropriate to perform trajectory-based learning?
- When is it appropriate to perform keyframe-based learning?



Trajectories and Keyframes for Kinesthetic Teaching: A Human-Robot Interaction Perspective



Sample demonstrations of the letter P in 2D

## **Trajectory Conversion**









Continuous trajectories in 2D

Data converted to keyframes

Clustering of keyframes and the sequential pose distributions Learned model trajectory

# Trajectory Conversion: Forward-Inverse Relaxation Model



Figure 3: An algorithm for extracting via-points.

- Fifth order splines used between positions to minimize jerk, using position, velocity, and acceleration per keyframe to compute the spline unknowns.
- Keyframes assume zero velocity/acceleration per point
- Trajectory demonstrations use the means from cluster centers.

A Computational Model for Cursive Handwriting Based on the Minimization Principle – Wada et al.

### Aligning Multiple Demonstrations



#### Implementation on PR2



#### Non-monolithic Task Representations

Can we learn and generalize multi-step tasks?

- Supports "life-long learning"
- Avoids dependency on isolated skill learning
  - Expensive to require human attention and demonstration
- Automatic segmentation allows for better skill transfer

Can we impart more complicated feature spaces into our skill representations without sacrificing usability?

## Skill Learning Wishlist

Recognize repeated instances of skills and generalize them to new settings.

Segment data without a priori knowledge of task structure.

Identify broad, general classes of skills (eg., manipulations, gestures, goal-based actions.)

Skill policies should have a flexible encoding such that they can be improved over time.

#### Learning and Generalization of Complex Tasks from Unstructured Demonstrations

[IROS 2015 / IJRR]

Scott Niekum, Sarah Osentoski, George Konidaris, Andrew G. Barto

Learning and Generalization of Complex Tasks from Unstructured Demonstrations

- Model-free skill segmentation
  - Using Bayesian nonparametric techniques
- Rapid policy learning
  - Learning from Demonstration accelerates skill acquisition
- Activity recognition without task priors
  - Using a Beta-Process Autoregressive Hidden Markov Model
- Flexible skill encoding
  - Dynamic Movement Primitives

# Task Learning Pipeline



- Task and skill representation created simultaneously from a continuous demonstration
- Recognizes re-used skills

### Preprocessing/Segmentation

• Segmentation is performed using a BP-AR-HMM [1]



### Skill Learning

- Skills are modeled as Dynamic Movement Primitives
  - Linear point attractor modulated by a nonlinear (learned) function
  - Uses end effector positions + quaternions for gripper rotation

$$f_{\text{target}}(s) = \frac{-K(g - x(s)) + D\dot{x}(s) + \tau \ddot{x}(s)}{g - x_0}$$



### **Evaluation: Testing Segmentation**



- Trained on demonstrations of the top task
- Tested on demonstrations of the bottom task



(a) Starting pose





(b) Reaches toward red (c) Picks up red block block (red block frame) (red block frame)







(e) Places red blockon green block(green block frame)



(f) Returns to home position (torso frame)

Autonomously segmented skills and associated frames

#### Failure Modes

- Symbolic Failure
  - Occurs when objects in the task description cannot be resolved (e.g., are missing from the environment)
  - Remedied through suggestion of substitutions
  - Possible corrections can be accepted or rejected due to pragmatic or preferential reasons.
    - Allows propositions for which the robot does not have an object model
- Execution Failure
  - Occurs when symbolic substitutions are accepted without a model sufficient for interaction
  - Occurs when outside the known policy region of a skill

#### **Application: Interactive Corrections**

#### INCREMENTAL SEMANTICALLY GROUNDED LEARNING FROM DEMONSTRATION

#### Abstraction is essential for solving complex problems



#### Not all robots operate in isolation









Autonomously Constructing Hierarchical Task Networks for Planning and Human-Robot Collaboration

[ICRA 15]

Bradley Hayes and Brian Scassellati

## Hierarchical Task Networks

#### **Benefits**

- Defines macro actions as compositions of primitive operators
- Provides a detailed problem factorization
- Operators are defined as (task, preconditions, effects)
- Facilitates look-ahead for increased execution flexibility (least-commitment planning)

#### **Challenges**

- Precise specifications of preconditions and effects can be difficult to specify
- Typically defined manually





#### Constructing Task Networks from Demonstrations

1. Extract task subgoals using min-cut

2. Convert task graph to subgoal graph

3. Apply a series of contraction operators to the subgoal graph

4. Create macro actions out of totally and partially ordered sub-plans at each iteration of contraction

## Extracting Sub-Goals: Intuition – Bottleneck Recognition





Problem Domain

State Frequency Map

[Q-Cut - Dynamic Discovery of Sub-Goals in Reinforcement Learning. Menache et al. 2002]

#### Application Domain - IKEA furniture





# Hierarchical Task Structure IKEA Chair



# Hierarchical Task Structure IKEA Chair



#### SMDP of "Attach Front Frame" Subtask





#### SMDP of "Attach Front Frame" Subtask



#### SMDP-Conjugate of "Attach Front Frame" Subtask



### **Building Hierarchical Structure**



Building a constraint-based hierarchy Goal:

Exploit existing structure to find logical groupings of sub-tasks

Task Graph to HTN Transform Input: SMDP-like Graph G Output: HTN H 1  $H \leftarrow \text{Conjugate-Graph-Transform}(G)$ 2 while  $|H_{vertices}| > 1$  do  $h\_size \leftarrow |H_{vertices}|$ 3 **foreach** maximal clique  $c = \{V, E\} \in H$  do 4 5 Compact  $\{V \in c\}$  into single metanode m 6 **foreach** maximal chain  $c = \{V, E\} \in H$  do 7 Compact  $\{V \in c\}$  into single metanode m 8 if  $h_{size} == |H_{vertices}|$  then break;

Step 0: Start Algorithm




State Action Any edges inbound to a clique member must have identical inbound edges to *all* clique members.

Any edge outbound from a clique member to an external vertex must have identical outbound edges from *all* clique members to the same target.



State Action Any edges inbound to a clique member must have identical inbound edges to *all* clique members.

Any edge outbound from a clique member to an external vertex must have identical outbound edges from *all* clique members to the same target.



Any edges inbound to a clique member must have identical inbound edges to *all* clique members.

Any edge outbound from a clique member to an external vertex must have identical outbound edges from *all* clique members to the same target.



State Action Any edges inbound to a clique member must have identical inbound edges to *all* clique members.

Any edge outbound from a clique member to an external vertex must have identical outbound edges from *all* clique members to the same target.

# Chains



Any edges inbound to a chain must only connect to the chain's starting vertex.

All internal nodes must have in and out degree 1.

Any edges outbound from the chain must only originate from the chain's terminating vertex.















#### Context-sensitive Supportive Behavior Policies



Interpretable Models for Fast Activity Recognition and Anomaly Explanation During Collaborative Robotics Tasks

> [ICRA 17] Bradley Hayes and Julie Shah

#### Collaborative robots need to recognize human activities

- Nearly all collaboration models depend on some form of activity recognition
- Collaboration imposes real-time constraints on classifier performance and tolerance to partial trajectories



Related Work

Fast target prediction of human reaching motion for cooperative human-robot manipulation tasks using time series classification

(Perez D'Arpino ICRA15)



## Related Work

Fast target prediction of human reaching motion for cooperative human-robot manipulation tasks using time series classification

(Perez D'Arpino ICRA15)



Time

Related Work

Fast target prediction of human reaching motion for cooperative human-robot manipulation tasks using time series classification

(Perez D'Arpino ICRA15)





Model each timestep with mean-centered Gaussian

Related Work

Fast target prediction of human reaching motion for cooperative human-robot manipulation tasks using time series classification

(Perez D'Arpino ICRA15)



### Common Activity Classifier Pipeline



- P. Koniusz, A. Cherian, and F. Porikli, "Tensor representations via kernel linearization for action recognition from 3d skeletons."
- Gori, J. Aggarwal, L. Matthies, and M. Ryoo, "Multitype activity recognition in robot-centric scenarios,"
- E. Cippitelli, S. Gasparrini, E. Gambi, and S. Spinsante, "A human activity recognition system using skeleton data from rgbd sensors."
- L. Xia, C. Chen, and J. Aggarwal, "View invariant human action recognition using histograms of 3d joints."

#### New Activity Recognition Approaches?

#### End-to-end Network

#### End-to-end Network

# In real deployments, humans need to be able to understand robot decisions



# In real deployments, humans need to be able to understand robot decisions

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG

Key Insight:

Take concepts from successful CNN/RNN classifiers and apply them to more transparent methods



Rapid Activity Prediction Through Object-oriented Regression (RAPTOR)

# A highly parallel ensemble classifier that is resilient to temporal variations















Displacement



Dictionary that maps IDs to sets of column indices E.g., {"Hands": [0,1,2,5,6,7]}

1	1.	0.04	-0.67	-0.4	-0.54	-0.74	-0.22	-0.75	-0.56
- (	0.04	1.	0.45	0.41	-0.03	-0.4	-0.44	-0.28	0.16
	-0.67	0.45	1.	0.39	0.49	0.2	-0.16	0.15	0.35
	-0.4	0.41	0.39	1.	0.06	0.2	0.11	0.13	0.38
	-0.54	-0.03	0.49	0.06	1.	0.36	0.02	0.16	0.39
	-0.74	-0.4	0.2	0.2	0.36	1.	0.37	0.57	0.19
	-0.22	-0.44	-0.16	0.11	0.02	0.37	1.	0.11	-0.25
- 1	-0.75	-0.98	0.15	0.13	0.16	0.57	0.11	1	0.57 L



Displacement



Within each temporal segment:

 Isolate columns of each demonstration trajectory according to (pre-defined) object map

( 1.	0.04	-0.67	-0.7	74 -0.22	-0.75
0.04	1.	0.45	-0.4	4 - 0.44	-0.28
-0.67	0.45	1.	0.2	-0.16	0.15
-0.4	0.41	0.39	0.2	0.11	0.13
-0.54	-0.03	0.49	0.36	0.02	0.16
-0.74	-0.4	0.2	1.	0.37	0.57
-0.22	-0.44	-0.16	0.37	1.	0.11
-0.75	-0.98	0.15	0.57	0.11	1

Create local model for each object



Displacement

Within each temporal-object segment:

- Ignore temporal information for each data point
- Treat as general pattern recognition problem
- Model the resulting distribution using a GMM

Result: An activity classifier ensemble across objects and time!



Need to find the most discriminative Object GMMs per time segment





Need to find the most discriminative Object GMMs per time segment



Need to find the most discriminative Object GMMs per time segment


# Activity Model Training Pipeline

- Choose top-N most discriminative features from the Random Forest classifier
- Weight each GMM proportional to its discriminative power



# Activity Model Training Pipeline

- Choose top-N most discriminative object-based classifiers
- Weight each object proportionally to its discriminative power



Result: Trained Highly Parallel Ensemble Learner with Temporal/Object-specific sensitivity



## Results: Three Datasets

- UTKinect publicly available benchmark (Kinect Joints)
- Dynamic Actor Industrial Manufacturing Task (Joint positions)
- Static Actor Industrial Manufacturing Task (Joint positions)







UTKinect

**Automotive Final Assembly** 

**Sealant Application** 

## Recognition Results: UTKinect-Action3D



Real-time UTKinect Activity Recognition Accuracy					
Classifier	Accuracy				
Slama et al. (2015) [21]	88.5%				
Chrungoo et al. (2014) [18]	89.45%				
Xia et al. (2012) [11]	90.9%				
Wang et al. (2015) [24]	90.9%				
Devanne et al. (2013) [20]	91.5%				
RAPTOR (proposed method)	92.1%				

pull	0.95	0	0	0	0	0	0	0	0.053	0 -
walk	0	1	0	0	0	0	0	0	0	0
push	- 0	0	0.68	0	0	0	0	0	0.32	0 -
pickUp	- 0	0.053	0	0.95	0	0	0	0	0	0 -
waveHands	0	0	0	0	1	0	0	0	0	0
carry	- 0	0.17	0	0	0	0.83	0	0	0	0 -
clapHands ·	- 0	0	0	0	0	0	0.95	0	0.053	0 -
standUp	- 0	0	0	0.053	0	0	0	0.95	0	0 -
throw	- 0	0	0	0	0	0	0.053	0	0.95	0 -
sitDown	- 0	0	0	0.053	0	0	0	0	0	0.95
	lInd	walk	- ysnd	pickUp -	sveHands	carry -	lapHands -	standUp -	throw	sitDown

## **Results: Online Prediction**



<b>RAPTOR Online Activity Prediction Accuracy</b>								
Dataset	25%	50%	75%	100%				
UTKinect	79.4%	83.1%	84.7%	92.1%				
Static-Reach	69.7%	77.2%	93.8%	97.5%				
Dynamic-AutoFA	91.7%	88.1%	90.5%	92.0%				

## Interpretability: Explaining Classifications

Key Insight:

- Apply outlier detection methods across internal activity classifiers
- Use outliers or lack thereof to explain issues across **time** and **objects**

Asking a "carry" classifier about a "walk" trajectory:

"In the **middle and end** of the trajectory, the **left hand and right hand** features were very poorly matched to my template."





Real-time Activity Segmentation and Classification



## Classification vs. Segmentation



What are the right intervals? Which intervals should get labels? Which labels should be where?



## A Naïve Changepoint Detection Approach

ScenarioDuration: 2700 frames1.5 minutes of dataClassifiers: 11– Avg run-time of 0.2s each

## IDEA: Run every activity classifier over every possible segment

- Given n frames:
  - For every interval q in the range [0, n]:
    - Evaluate each classifier on *q* —
  - Sort results by likelihood
  - Assign class labels to uncovered intervals from highest likelihood classifications until no unlabeled frames remain
- Return timeline (list of intervals)

2700<sup>2</sup> \* 0.2 = 1458000sec ~16.88 days

Classifiers must be ideal (sensitive to trajectory length, non-overlapping, comparable tolerance to noise, etc.)



particle\_maps[] - Sorted (MAP, particle) tuples for each timestep

- At each time step *t*:
  - Create new particles for all eligible classes
    - start\_time = t minimum\_class\_duration
    - prev\_interval = particle with highest MAP estimate in best[start\_time]
  - Evaluate existing particles' likelihoods over the interval [*p.start\_time, t*] and store as (likelihood, p) tuples in *particle\_maps*
  - Terminate stale particles



particle\_maps[] - Sorted (MAP, particle) tuples for each timestep

- At each time step *t*:
  - Create new particles for all eligible classes
    - *start\_time = t minimum\_class\_duration*
    - prev\_interval = particle with highest MAP estimate in best[start\_time]
  - Evaluate existing particles' likelihoods over the interval [*p.start\_time, t*] and store as (likelihood, p) tuples in *particle\_maps*
  - Terminate stale particles



particle\_maps[] - Sorted (MAP, particle) tuples for each timestep

- At each time step *t*:
  - Create new particles for all eligible classes
    - *start\_time = t minimum\_class\_duration*
    - prev\_interval = particle with highest MAP estimate in best[start\_time]
  - Evaluate existing particles' likelihoods over the interval [*p.start\_time, t*] and store as (likelihood, p) tuples in *particle\_maps*
  - Terminate stale particles



- At each time step *t*:
  - Create new particles for all eligible classes
    - *start\_time* = *t minimum\_class\_duration*
    - prev\_interval = particle with highest MAP estimate in best[start\_time]
  - Evaluate existing particles' likelihoods over the interval [*p.start\_time, t*] and store as (likelihood, p) tuples in *particle\_maps*
  - Terminate stale particles



- Set *f* = final frame index
- While *f* > 0 and *particle\_maps[f]* != None:
  - Take best (MAP, particle) at particle\_maps index f
  - Annotate segment [*shape\_start, f*] with *shape\_class*
  - Set *f* = *particle.start\_time*



- Set *f* = final frame index
- While *f* > 0 and *particle\_maps[f]* != None:
  - Take best (MAP, particle) at particle\_maps index f
  - Annotate segment [*shape\_start, f*] with *shape\_class*
  - Set *f* = *particle.start\_time*



- Set *f* = final frame index
- While *f* > 0 and *particle\_maps[f]* != None:
  - Take best (MAP, particle) at particle\_maps index f
  - Annotate segment [*shape\_start, f*] with *shape\_class*
  - Set *f* = *particle.start\_time*



- Set *f* = final frame index
- While *f* > 0 and *particle\_maps[f]* != None:
  - Take best (MAP, particle) at particle\_maps index f
  - Annotate segment [*shape\_start, f*] with *shape\_class*
  - Set *f* = *particle.start\_time*



- Set *f* = final frame index
- While *f* > 0 and *particle\_maps[f]* != None:
  - Take best (MAP, particle) at particle\_maps index f
  - Annotate segment [*shape\_start, f*] with *shape\_class*
  - Set *f* = *particle.start\_time*

## Associating Robot Behaviors with Task States



## Motion models in collaborative settings

[Hayes and Scassellati ICDL13]



# At a high level, *social force* is a projection of an agent's physical space occupation via its anticipated travel path

Social force carries different meanings depending on the task and environmental contexts in which it is applied.

# Motion models in collaborative settings

## Field treatment dictates robot's role



## Social Force in Human-Robot Teaming



# Take a break (10 minutes)

Policy Shaping:

- Stand up
- Get Caffeine
- Go stand outside for a few minutes
- (Talk to someone about how IML relates to your work)

## Cooperative Inverse Reinforcement Learning

[NIPS 2016]

Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, Stuart Russell

Inverse Reinforcement Learning Can Break Down in Team Scenarios!

- Traditional IRL is optimal if the reference demonstrations are "Expert" demonstrations.
  - ...but execution happens in isolation!
- Expert demonstrations are not always the most effective teaching strategy:
  - Sometimes better to learn the landscape of the problem than to see a optimal demonstrations
- Properly crafted 'imperfect' demonstrations can better communicate information about the objective

## The Shutdown Problem







Slide credit: Dylan Hadfield-Menell, "Cooperative Inverse Reinforcement Learning", CITRIS Workshop on Algorithmic Human-Robot Interaction. INRIA 2016.

## The Shutdown Problem



$$P(sd)R_{sd} + (1 - P(sd))R$$

$$R_{sd}$$

Slide credit: Dylan Hadfield-Menell, "Cooperative Inverse Reinforcement Learning", CITRIS Workshop on Algorithmic Human-Robot Interaction. INRIA 2016.

## Issues in Inverse Reinforcement Learning

- Uncertainty about task objectives is essential for cooperative behaviors
- IRL Pitfalls:
  - Don't want to just imitate the demonstrator
  - Assumes the demonstrator is 'unaware' of being observed
  - Action selection is independent of reward uncertainty
- Without modeling reward uncertainty, robot gets narrow view of environment dynamics and reward
- From Ramachandran and Amir's "Bayesian Inverse Reinforcement Learning":
  - The optimal policy for an MDP with a distribution over reward functions R ~ P(R) is one that maximizes reward according to the expectation of R.

## Proposal: Robot Plays Cooperative Game

- Cooperative Inverse Reinforcement Learning
  - [Hadfield-Menell et al. NIPS 2016]
- Two players:



- Both players maximize a shared reward function, but only H observes the actual reward signal; R only knows a prior distribution on reward functions
  - ${f R}$  learns the reward parameters by observing  ${f H}$

**Cooperative Inverse Reinforcement Learning** 



Cooperative Inverse Reinforcement Learning

$$\langle \mathcal{S}, \{\mathcal{A}^{\mathbf{H}}, \mathcal{A}^{\mathbf{R}}\}, T, \{R, \Theta, P_0\}, \gamma \rangle$$

- t=-1  $\theta \sim P_0(\theta)$
- t=0  $\mathbf{H}$  observes  $\theta$
- For t = 0, ...
  - $\mathbf{H}$  and  $\mathbf{R}$  observe  $S_{\mathbf{T}}$
  - $\mathbf{H}$  and  $\mathbf{R}$  select  $a_t^{\mathbf{H}}$  and  $a_t^{\mathbf{R}}$  respectively
  - New state  $S_{t+1}$  is sampled from  $T(\cdot | s_t, a_t^H, a_t^R)$
  - Both observe each other's actions and collect reward  $R(s_t; heta)$

## Cooperative Inverse Reinforcement Learning



## **CIRL** Properties

- The distribution over state sequences is determined by a pair of policies: (  $\pi^{H}$ ,  $\pi^{R}$  )
- An 'optimal' policy pair maximizes the discounted sum of rewards
- In general, policies may depend on the entire observation histories
  - The history of states and actions for both actors includes the reward parameter for the human
  - [Hadfield-Menell '16] There exists an optimal policy pair that only depends on the current state and the robot's belief

## Incentives for Instructive Demonstrations

- Reduces the robot's expected regret
- Reduces the KL Divergence of trajectory distributions
- Reduces reward errors

Further reading:

Game-Theoretic Modeling of Human Adaptation in Human-Robot Collaboration by Nikolaidis et al.

HRI 2017

Extends CIRL, providing a model of human partial adaptation to a robot collaborator without adopting the robot's policy as their own.





## Effective Robot Teammate Behaviors for Supporting Sequential Manipulation Tasks

## [IROS 2015] Bradley Hayes and Brian Scassellati
### Can we do better than LfD for Skill Acquisition?



Human figures out *how* and *when* the robot can be helpful

Quickly enables useful, helpful actions.

Does not scale with task count! Requires human expert

Robot figures out *how* and *when* it can be helpful

- Allows for novel behaviors to be discovered
- Enables deeper task comprehension and action understanding

### Autonomously Generating Supportive Behaviors: A Task and Motion Planning Approach



#### Go to object bx GOTOB(bx) Preconditions: TYPE(bx,OBJECT),(3rx)[INROOM(bx,rx) & INROOM(ROBOT,rx)] Deletions: AT(ROBOT,\$1,\$2), NEXTTO(ROBOT,\$1) \*NEXTTO(ROBOT,bx) Additions Go to door dx GOTOD(dx) Preconditions: TYPE(dx,DOOR),(3rx)(3ry)[INROOM(ROBOT,rx) & CONNECTS(dx,rx,ry)] Deletions: AT(ROBOT,\$1,\$2), NEXTTO(ROBOT,\$1) Additions \*NEXTTO(ROBOT,dx) Go to coordinate location (x,y). GOTOL(x,y) Preconditions: (3rx)[INRODM(ROBOT,rx) & LOCINROOM(x,y,rx)] Deletions: AT(ROBOT,\$1,\$2), NEXTTO(ROBOT,\$1) Additions: \*AT(ROBOT,x,y) Go through door dx into room rx. GOTHRUDR(dx,rx) Preconditions: TYPE(dx,DOOR), STATUS(dx,OPEN), TYPE(rx,ROOM), NEXTTO(ROBOT,dx) (3rx)[INROOM(ROBOT,ry) & CONNECTS(dx,ry,rx)] Deletions AT(ROBOT,\$1,\$2), NEXTTO(ROBOT\$1), INROOM(ROBOT,\$1) Additions: \*INROOM(ROBOT,rx)



#### **Perspective Taking**

#### Symbolic planning

#### Motion planning

Autonomously Generated Supportive Behaviors

# Task and Motion Planning

The TAMP problem is represented by the tuple: {A, O, C, s<sub>0</sub>, s<sub>G</sub>}

A is a set of (lead) agents

**O** is a set of operators (unparameterized motor primitives)

**C** is a capabilities mapping function between agents and operators

 $\mathbf{s}_0$  is the set of predicates *precisely* specifying the start state

 $\mathbf{s}_{\mathbf{G}}$  is the set of predicates specifying the goal state

# Supportive Behavior TAMP

The SB-TAMP problem is represented as the tuple: { T,  $\Pi_T$ ,  $a_s$ ,  $C_s$ ,  $s_c$ , P }

- **T** is a TAMP problem
- Π<sub>T</sub> is a set of symbolic plans for T
- **a**<sub>s</sub> is a supportive agent
- C<sub>s</sub> is a mapping function indicating operators from T usable by a<sub>s</sub>
- **s**<sub>c</sub> is the current environment state
- **P** is a set of partially or fully specified predicates describing prohibited environmental effects for support actions

## Supportive Behavior Pipeline: Intuition



- Propose alternative environments

   Change one thing about the environment
- Evaluate if they facilitate the leader's task/motion planning

   Simulate policy execution(s) from leader's perspective
- 3. Compute cost of creating target environment

- Simulate support agent's plan execution

- 4. Choose environment that maximizes [benefit cost]
  - Execute supportive behavior plan

# Supportive Behavior Pipeline



# Supportive Behavior Pipeline



# Plan Evaluation

Choose the support policy ( $\xi \in \Xi$ ) that minimizes the expected execution cost of the leader's policy ( $\pi \in \Pi$ ) to solve the TAMP problem **T** from the current state ( $s_c$ )

- Cost estimate must account for
  - Resource conflicts (shared utilization/demand)
  - Spatial constraints (support agent's avoidance of lead)

$$\min_{\xi \in \Xi} \sum_{\pi \in \Pi_T} w_\pi * \operatorname{cost} (T, \pi, \xi, s_c, \gamma)$$

# Plan Evaluation

Choose the support policy  $(\xi \in \Xi)$  that minimizes the expected execution cost of the leader's policy  $(\pi \in \Pi)$  to solve the TAMP problem **T** from the current state  $(s_c)$ 



# Weighting functions: Uniform, Greedy

$$w_{\pi}$$
 = 1

Consider all known solutions equivalently likely and important

$$w_{\pi} = \begin{cases} 1 ; & \operatorname{duration}(T, \pi, \emptyset, s_0, f(x) = 1) = \underset{\operatorname{duration}}{\overset{\operatorname{Min}}{}} \\ 0 ; & \operatorname{otherwise} \end{cases}$$

Only the best-known solution is worth planning against

# Weighting functions: Uniform



# Weighting functions: Optimality-Proportional

$$w_{\pi} = \left(\frac{\min_{\pi \in \Pi_{T}} \operatorname{duration}(T, \pi, \emptyset, s_{0}, f(x) = 1)}{\operatorname{duration}(T, \pi, \emptyset, s_{0}, f(x) = 1)}\right)^{\mathsf{p}}$$

Weight plans proportional to similarity vs. the best-known solution



# Weighting functions: Optimality-Proportional



# Weighting functions: Error Mitigation

$$w_{\pi} = \begin{cases} f(\pi) & ; \text{ duration}(T, \pi, \emptyset, s_0, f(x) = 1) \leq \epsilon \\ -\alpha w_{\pi} & ; \text{ otherwise} \end{cases}$$

Plans more optimal than some cutoff  $\varepsilon$  are treated normally, per *f*.

Suboptimal plans are **negatively weighted**, encouraging active mitigation behavior from the supportive robot.

 $\alpha < \frac{1}{\max_{\pi} w_{\pi}}$  is a normalization term to avoid harm due to plan overlap

# Weighting functions: Error Mitigation



# Effect of Supportive Behaviors

Task Completion Time



## Improving Robot Controller Transparency Through Autonomous Policy Explanation

[HRI 2017] Bradley Hayes and Julie Shah

## Shared Expectations are Critical for Teamwork

In close human-robot collaboration...

- Human must be able to plan around expected robot behaviors
- Understanding failure modes and policies are central to ensuring safe interaction and **managing risk**

Fluent teaming **requires** communication...

- When there's no prior knowledge
- When expectations are violated
- When there is joint action







When will you stop helping me pour the molten aluminum?













I will terminate **assist aluminum pouring** when the world is in state:

$\langle \rangle$	$\langle \rangle$	$\langle \rangle$
12.4827	15	12.4827
5.12893	7.125	8.51422
1.12419	1.12419	1.12419
0	0	0
0	0	1
1	1	0
3.62242	-8.1219	3.62242
-40.241	-40	-40.241
	,,	

### State space is too obscure to directly articulate



# Motivation







## Motivation

### Why did the robot not inspect the orange gear?



## Motivation

### How do we diagnose and repair this fault?



# State of the Art



int \*detect\_gear = &INPUT1; int \*gear\_x = &INPUT2;

if (\*detect\_gear == 1 && \*gear\_x <= 10 && \*gear\_x >= 8) {
 pick\_gear(gear\_x);

???



# Establishing Shared Expectations





Role-based Feedback [St. Clair et al. 2016] Legible Motion [Dragan et al. 2013]



State Disambiguation [Wang et al. 2016]

	BREAK	Idle	NEGOTIATE	Sell	INNERTALK	WATCH	GUARD	EQUIP
ANNY	( 0	0	0	1	1	1	0	0
BENNY	0	0	0	1	1	1	0	0
CANNY	1	0	0	0	0	0	0	1
DANNY	0	0	1	1	1	0	0	0
ERNY	0	1	0	0	0	0	1	0
FRENNY	\ 1	0	0	0	0	0	0	1

Coordination Graphs [Kalech 2010]



Hierarchical Task Models [Hayes et al. 2016]





Policy Dictation [Johnson et al. 2006] Collaborative Planning [Milliez et al. 2016]

#### Short Term

#### Long Term

# Natural Interaction

### **Reasonable question:**

"Why didn't you inspect the gear?"

### Interpretable answer:

"My camera didn't see a gear. I inspect the gear when it is less than 0.3m from the conveyor belt center and it has been placed by the gantry."





# Making Control Systems More Interpretable

### Approach:

- 1. Attach a smart debugger to monitor controller execution
- 2. Build a graphical model from observations
- 3. Use specialized algorithms to map queries to state regions
- 4. Collect relevant state region attributes
- 5. Minimally summarize relevant state regions with attributes
- 6. Communicate query response

	Model Building
ns	
	Query Analysis
tes	
	Response Generation

## **Expectation Synchronization**



### Required

Policy model

Concept representations

Mapping from query to model

Mapping from model to response

# Policy Modeling

### Local, approximate behavioral models from observation

(Generate MDP from regular controller operation)



Deployment or simulation environment

**Behavioral MDP** 

Control policy

States are composed of internal variables and externally sensed information

Actions are parameterized function calls observed from the controller

Transitions are learned by observing resultant states from function calls

### Given



### Required

Policy model

Concept representations

Mapping from query to model

Mapping from model to response
# **Concept Representations**

**Concept library**: generic state classifiers mapped to semantic templates that identify whether a state fulfills a given criteria

Set of Boolean classifiers:

- Spatial concepts
- Domain-specific concepts
- Agent-specific concepts



camera powered

(e.g., "A is on top of B") (e.g., "Widget paint is drying")

State  $\rightarrow$  {True, False}



### Given



### Required

Policy model

**Concept representations** 

Mapping from query to model

Mapping from model to response

### Improving Control Policy Transparency

Three template questions for synchronizing expectations:

- When do you {*action*}?
- What will you do when
  {environmental conditions}?
- Why didn't you do {*action*}?



# **Relevant Question Templates**

### When will you do {action}?



Algorithm 2: Identify Dominant-action State Region

**Input:** Behavioral Model  $G = \{V, E\}$ , Target Action  $a_t$ **Output:** Set of target states  $S_{\pi^a}$ , Set of non-target states

1 
$$S_{\pi^a} \leftarrow \{\};$$

$$2 S_{\pi^* \backslash a} \leftarrow \{\};$$

3 foreach  $s \in V$  do

 $a \leftarrow \text{most frequent action executed from } s;$ 

5 if 
$$a == a_t$$
 then  $S_{\pi^a} \leftarrow S_{\pi^a} \cup s$ ;

6 else 
$$S_{\pi^* \setminus a} \leftarrow S_{\pi^* \setminus a} \cup s$$
;

7 return  $S_{\pi^a}, S_{\pi^* \setminus a}$ ;

# **Relevant Question Templates**

### Why didn't you do {action}?



Algorithm 3: Identify Behavioral Divergences **Input:** Behavioral Model  $G = \{V, E\}$ , Target Action  $a_t$ , Previous state  $s_p$ , Distance threshold  $D_{const}$ Output: Explanation of difference between current state and state region where  $a_t$  is performed, explanation of where  $a_t$  is performed locally. 1  $S_{\pi^a} \leftarrow \{\};$ 2  $S_{\pi^* \setminus a} \leftarrow \{\};$ 3 foreach  $D \in \{1, ..., D_{const}\}$  do foreach  $s \in \{v \in V \mid distance(v, s_p) \leq D\}$  do  $a \leftarrow \text{most frequent action executed from } s;$ 5 if  $a == a_t$  then  $S_{\pi^a} \leftarrow S_{\pi^a} \cup s$ ; else  $S_{\pi^* \setminus a} \leftarrow S_{\pi^* \setminus a} \cup s$ ; 6 7 s expected\_region  $\leftarrow$  describe $(G, S_{\pi^a}, S_{\pi^* \setminus a});$ 9 current region  $\leftarrow$  describe $(G, \{s_n\}, S_{\pi^a})$ ; 10 return diff(expected\_region, current\_region), expected region;

# **Relevant Question Templates**

### What will you do when {conditions}?

5

6

7

9



**Input:** Behavioral Model  $G = \{V, E\}$ , Concept Library C, State region description d, Max action threshold *cluster* max **Output:** Explanation of behavior in d, broken down by action and accompanying state region 1  $S \leftarrow dict();$ 2 descriptions  $\leftarrow$  dict(); 3 DNF description  $\leftarrow$  convert to DNF formula(d, C); 4 foreach  $s \in \{v \in V \mid test dnf(v, DNF description) is$ True } do  $S[\pi(s)] \leftarrow S[\pi(s)] \cup s;$ if |S| > cluster max then return too\_many\_actions\_error s foreach  $a \in S$  do descriptions[a]  $\leftarrow$  describe(S[a]); 10 return descriptions;

Algorithm 4: Characterize Situational Behavior

### Given



### Required

Policy model

Concept representations

Mapping from query to model

Mapping from model to response

### Language Mapping: Model to Response

# **Recall**: Concept library provides dictionary of classifiers that cover state regions



### Using Concepts to Describe State Regions

We perform state-to-language mapping by applying

a Boolean algebra over the space of concepts



This reduces concept selection to a set cover problem over state regions

Disjunctive normal form (DNF) formulae enable coverage over arbitrary geometric state space regions via **intersections** and **unions** of concepts

Templates provide a mapping from DNF  $\rightarrow$  natural language

# Query Response Process



# **Producing Efficient Summaries**

### Achieving the succinctness criterion is **NP-hard**.

(choosing the minimal set of concepts with the best state region coverage precision)



The same problems in succinctness are encountered during circuit minimization:

Prime implicants are concept clauses covering minterms (target states)

Can use Quine-McCluskey Algorithm to find minimization

Q-M doesn't scale well, but we can get approximate solutions using ESPRESSO, with appropriate sacrifices of optimality or precision.

### Given



### Required

Policy model

**Concept representations** 

Mapping from query to model

Mapping from model to response



### Result: Agents can explain their policies to collaborators



"I will get the gear when I am near a human-only zone and I do not carry a gear. I will move north when I am near a humanonly zone and I carry a gear."



"I move right when the cart is at the far left or when the cart is in the middle and the pole is falling right or when the cart is in the far right and the pole is stabilizing left."



"I didn't inspect the part because the stock feed signal is off. I inspect the part when the stock feed signal is on and I have detected a part and the part is within reach."

### Designing Interactions for Robot Active Learners

Maya Cakmak, Crystal Chao, Andrea L. Thomaz [TAMD 2010]

# Passive vs. Active Learning

- Active Learning embeds robots in a tightly coupled dyadic interaction
  - Improper handling of this interaction disengages the oracle!
  - Disengagement leads to poor information quality
- Difficult to balance learning accuracy and learning speed with interaction smoothness
- Active Learning is a tool for increasing learner transparency

## Experiment: Teaching Shape Composites

Concept	Concept	Examples	# of
Name	Representation		Instances
HOUSE	$shape_{top} = triangle \land \ color_{top} = pink \land \ shape_{bottom} = square$		16
SNOWMAN	$shape_{top} = circle \land$ $size_{top} = small \land$ $shape_{bottom} = circle$	<b>e</b>	28
ALIEN	$shape_{top} = circle \land color_{top} = green \land color_{bottom} = green$		10
ICE CREAM	$shape_{top} = circle \land$ $shape_{bottom} = triangle \land$ $color_{bottom} = yellow$		16



# Version Space Learning

Step	Example	Label	Version Space
1	<pink,triangle,large, yellow,square,large&gt;</pink,triangle,large, 	+	Most specific hypothesis: <pink,triangle,large yellow,square,large&gt; Most general hypotheses: <pink,*,*,*,*,*> &lt;*,triangle,*,*,*,*&gt;  &lt;*,*,*,*,*,large&gt;</pink,*,*,*,*,*></pink,triangle,large 
2	<pre><orange,triangle,large, yellow,square,large=""></orange,triangle,large,></pre>	-	Most specific hypothesis: <pink,triangle,large yellow,square,large&gt; Most general hypothesis: <pink,*,*,*,*,*></pink,*,*,*,*,*></pink,triangle,large 
3	<pink,triangle,small, yellow,square,small&gt;</pink,triangle,small, 	+	Most specific hypothesis: <pink,triangle,* yellow,square,*&gt; Most general hypothesis: <pink,*,*,*,*,*></pink,*,*,*,*,*></pink,triangle,* 

Example	Version Space Predictions	Total	Prediction
<pink,triangle,large yellow,square,small&gt;</pink,triangle,large 	< pink, triangle, * yellow, square, $* > \rightarrow +$ $\dots$ $< pink, *, *, *, *, * > \rightarrow +$	+: 8 -: 0	+
<pre><pink,circle,large yellow,circle,small=""></pink,circle,large></pre>	< pink, triangle, * yellow, square, $* > \rightarrow -$ $\dots$ $< pink, *, *, *, *, * > \rightarrow +$	+: 2 -: 6	-
<pink,triangle,small green,square,small&gt;</pink,triangle,small 	< pink, triangle, * yellow, square, $* > \rightarrow -$ $\dots$ $< pink, *, *, *, *, *> \rightarrow -$	+: 4 -: 4	?

Teacher Actions: { "This is a <concept>." | "This is not a <concept>." | "Is this a <concept>?" }

# Active Learning Gets Better Coverage

Mode	Final	Subjects who	Number of Steps
	$F_1$ -scores	achieved $ V  = 1$	to reach $ V  = 1$
SL	77.81%	6/24 (25%)	M = 10.67, SD = 4.59
AL	100.00%	24/24 (100%)	M = 8.29, SD = 2.29
MI	98.61%	23/24 (96%)	M = 8.35, SD = 1.56
AQ	97.50%	20/24 (83%)	M = 9.30, SD = 3.53

Humans have a tough time keeping track of their teaching progress, even for small instance spaces.

# Subjective Measures

Mode	Intelligence (7=best)	Enjoyability (7=best)
SL	M = 4.08, $SD = 1.64$	M = 4.92, SD = 1.79
AL	M = 5.54, SD = 1.53	M = 5.46, SD = 1.77
MI	M = 5.12, SD = 1.57	M = 5.46, SD = 1.50
AQ	M = 5.75, SD = 1.29	M = 6.04, $SD = 1.16$

Active Learning modes are perceived as both more intelligent and more enjoyable to interact with

People preferred control over triggering the robot's Active Learning mechanism.