# Acquiring Rich Models of Objects and Space Through Vision and Natural Language

## Matthew Walter

CS & AI Lab, MIT
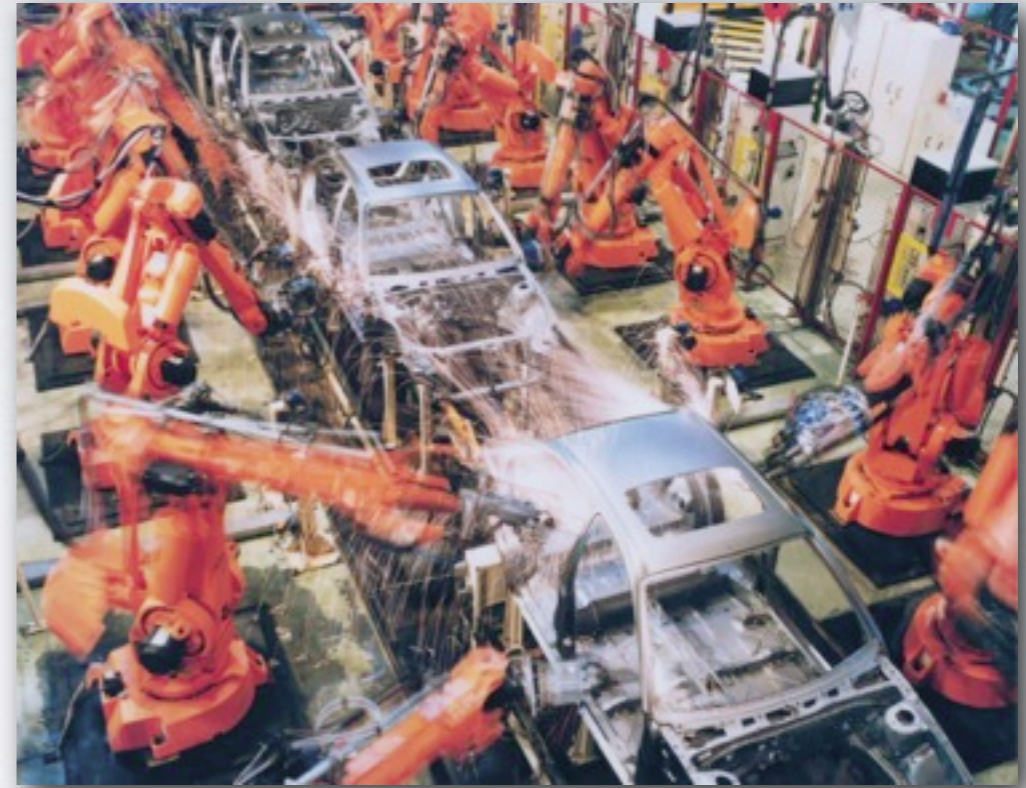
School of Computer Science
University of Massachusetts, Amherst

February 6, 2013

|

CSAIL

# Robots as Automated Agents

- Advances in:
  - Actuation
  - Planning
  - Control

- Focus:
  - Accuracy
  - Robustness



Courtesy: ABB

Matthew Walter
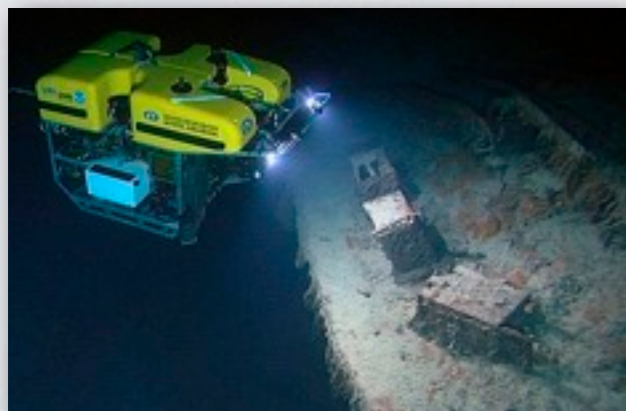
Tuesday, February 5, 13

# Robots as Our Surrogates

- Advances in:
  - Estimation
  - Navigation
  - Planning under uncertainty

- Focus:
  - Accuracy
  - Robustness



[JFR 2008]





RMS Titanic

[IJRR 2006]

Matthew Walter

# Robots as Our Partners



THE BOSTON GLOBE
**Business**
Science & Innovation
**Wheelchairs that listen**

William Li of MIT (left) and David Hatch at The Boston Home

PHOTO BY JOHN TLUMACKI/GLOBE STAFF



CSAIL ENVOY

Matthew Walter

Tuesday, February 5, 13

# Now: People Accommodating Robots



Courtesy: AeroVironment



Courtesy: US Army

Matthew Walter

Tuesday, February 5, 13

# Where We Need to Be

Matthew Walter

Tuesday, February 5, 13

# Where We Need to Be

Matthew Walter

Tuesday, February 5, 13

# Representational Divide

## A robot's view of the world is very different from our own

People

Robots

- Objects
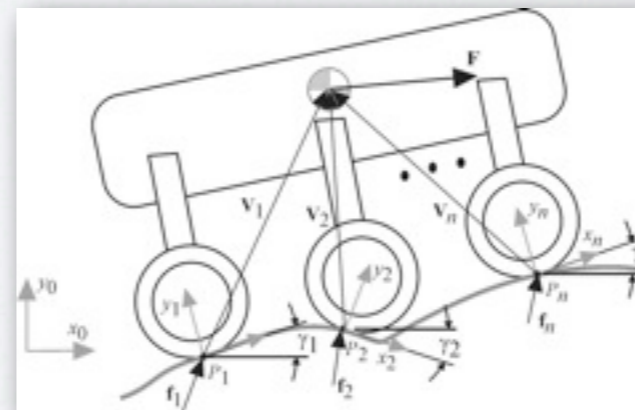
- Places

- Actions

- People

- Events


Images


Laser scans


Wheel torques


Figure 2 The unimate PUMA 562 robot arm
Joint angles

Matthew Walter

Tuesday, February 5, 13

# Vision: Shared Situational Awareness

Spatially extended, temporally persistent
model of the robot's surround

- Objects: Identity, properties, relations, actions

- Places: Function, identity, connectivity

- People: Locations, behavior, gestures

- Actions: Means of interacting with the world

Matthew Walter

Tuesday, February 5, 13

# Vision: Learning Shared Representations

- Reason over shared knowledge representations

- Acquire situational awareness as they interact with the world

- Learn opportunistically from humans

Matthew Walter

Tuesday, February 5, 13

Matthew Walter

Tuesday, February 5, 13

# Assistive Mobile Manipulation





Courtesy: University of Pittsburgh

- Material handling in unstructured environments

- Assisted living for the elderly & disabled

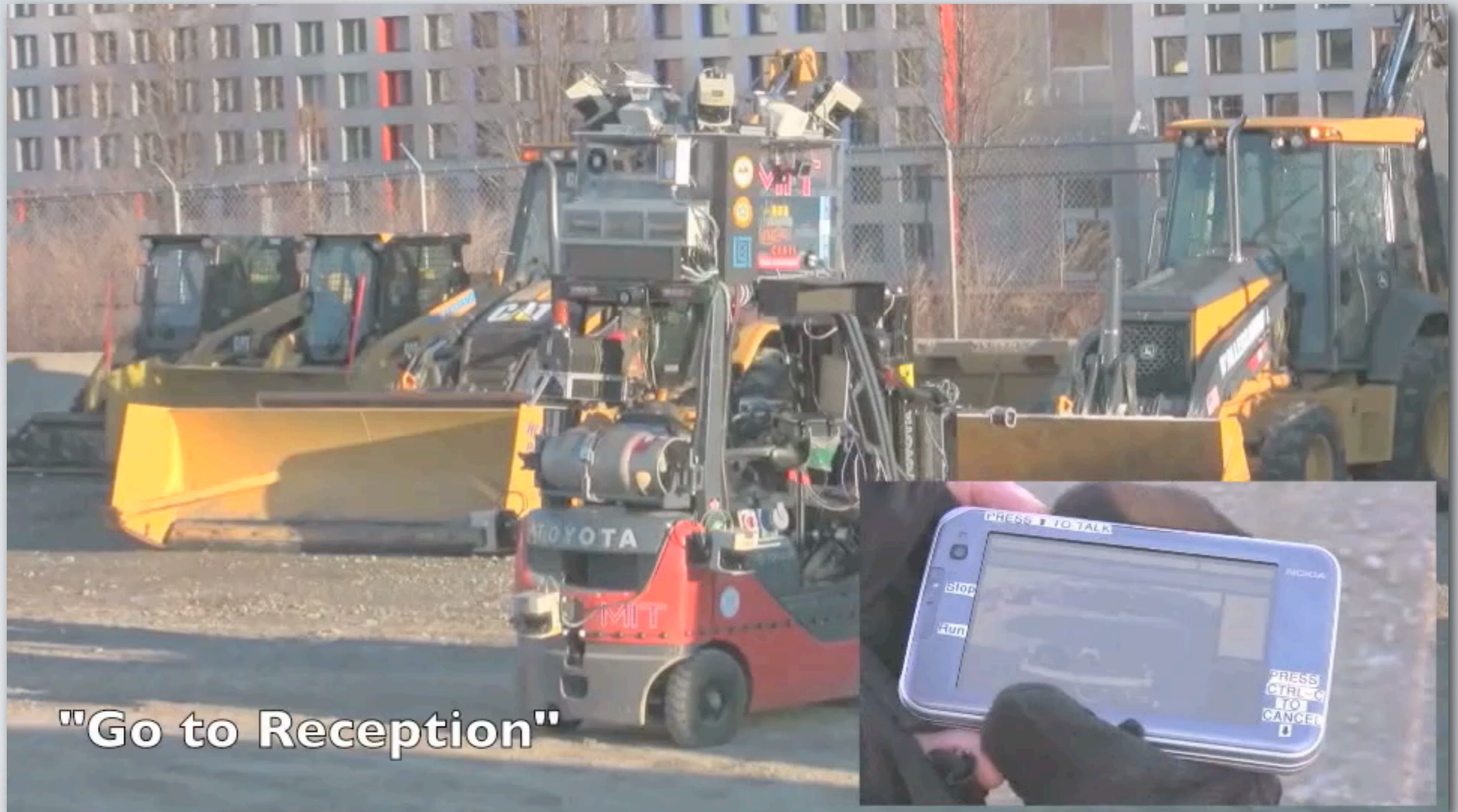Matthew Walter

Tuesday, February 5, 13

# Challenges for Mobile Manipulation





- Unprepared, dynamic environments
  - Coarse localization
  - Uneven terrain
  - Uncontrolled lighting

- Objects unknown a priori

- People everywhere

- Intuitive, human-centered control

Matthew Walter

Tuesday, February 5, 13

# Shared Autonomy



"Go to Reception"

Matthew Walter

Tuesday, February 5, 13

# Efficient Manipulation via Object Awareness



Put the pipes on the truck

Please pick up my keys

Courtesy: Kinova Robotics

Matthew Walter

Tuesday, February 5, 13

# Object Recognition is Hard!

- Usability requirements:
  - Persistent, reliable detection
  - Efficient object learning

- Challenges:
  - Variable lighting (outdoors)
  - Variable viewpoints
  - Unobserved object relocation
  - Coarse localization

Matthew Walter

Tuesday, February 5, 13

# Object Recognition is Hard!

- Usability requirements:
  - Persistent, reliable detection
  - Efficient object learning

- Challenges:
  - Variable lighting (outdoors)
  - Variable viewpoints
  - Unobserved object relocation
  - Coarse localization

Matthew Walter

Tuesday, February 5, 13

# Object Recognition is Hard!

- Usability requirements:
  - Persistent, reliable detection
  - Efficient object learning

- Challenges:
  - Variable lighting (outdoors)
  - Variable viewpoints
  - Unobserved object relocation
  - Coarse localization

Matthew Walter

# Object Instance Recognition

| | Object category detection [1] | Visual tracking [2] | **This work** [3] |
|---|---|---|---|
| Train from one example | | ✓ | ✓ |
| Train online | | ✓ | ✓ |
| Persistence (hours/days) | | | ✓ |
| Category recognition | ✓ | | |
| Instance recognition | | ✓ | ✓ |
| Real-time performance | | ✓ | ✓ |

[1] Nistér'06, Hoiem'07, Savarese'07

[2] Collins'05, Grabner'08, Kalal'09

[3] CVPRW'10, ISER'10, IJRR'12

Matthew Walter

Tuesday, February 5, 13

# One-shot Appearance Learning

- Key ideas for usable object reacquisition
  - Detect *instances* of the objects used in practice
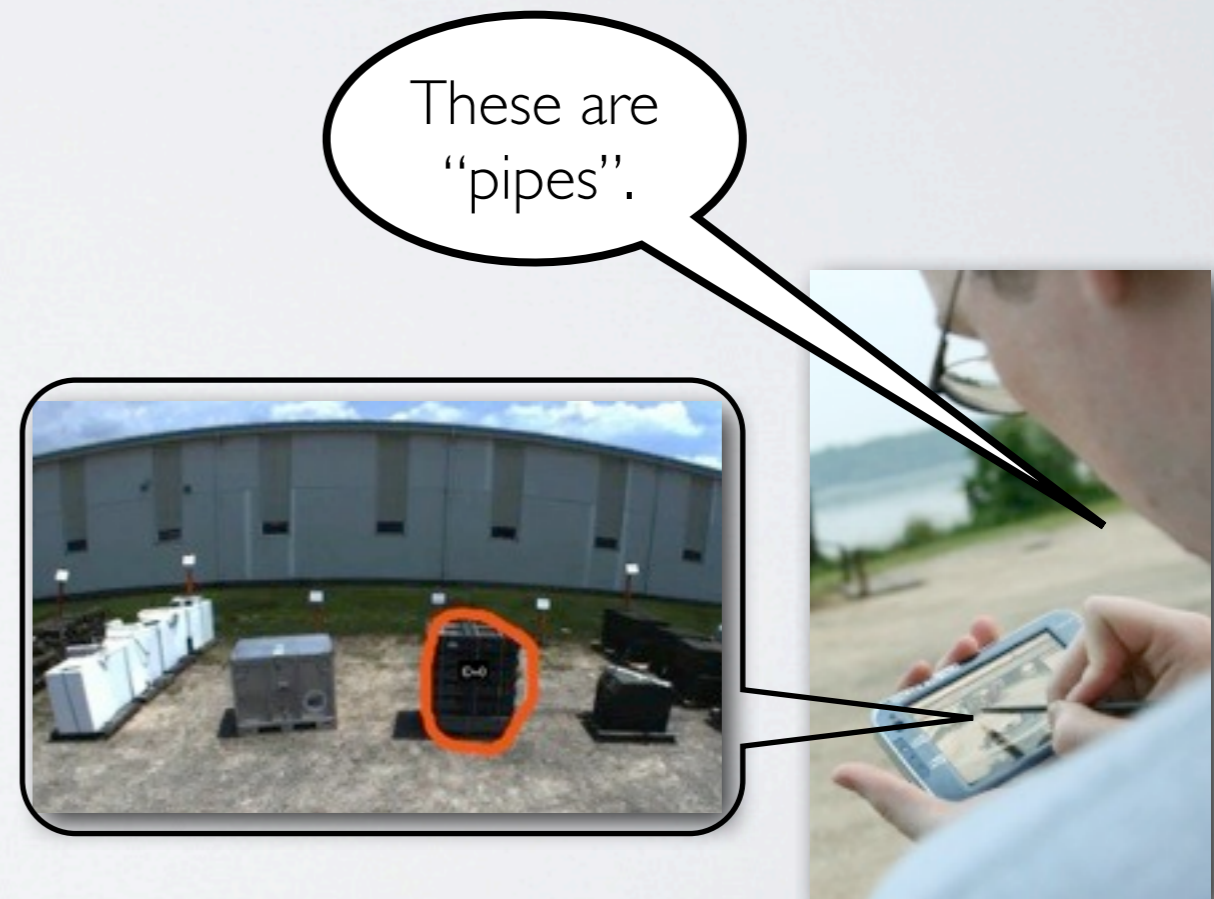  - Take advantage of the robot's mobility for learning

Matthew Walter

Tuesday, February 5, 13

# One-shot Appearance Learning

- Key ideas for usable object reacquisition
  - Detect *instances* of the objects used in practice
  - Take advantage of the robot's mobility for learning

Matthew Walter

Tuesday, February 5, 13

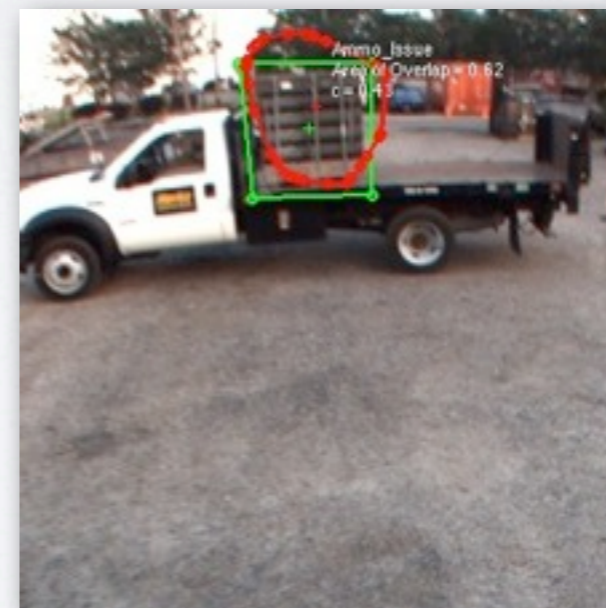# One-shot Appearance Learning

- User provides a single example of the object's identity (name & segmentation)

- System bootstraps on user's example to build an appearance model online

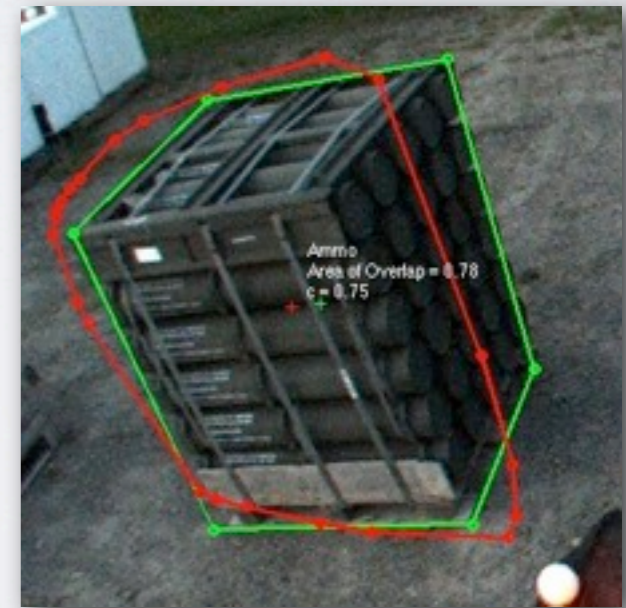- System takes advantage of robot's motion to opportunistically capture appearance variations



These are "pipes".

[ISER 2010; IJRR 2012]

Matthew Walter

Tuesday, February 5, 13

# One-shot Appearance Learning

- User provides a single example of the object's identity (name & segmentation)

- System bootstraps on user's example to build an appearance model online

- System takes advantage of robot's motion to opportunistically capture appearance variations



Illumination                 Aspect                 Relocation                 Scale
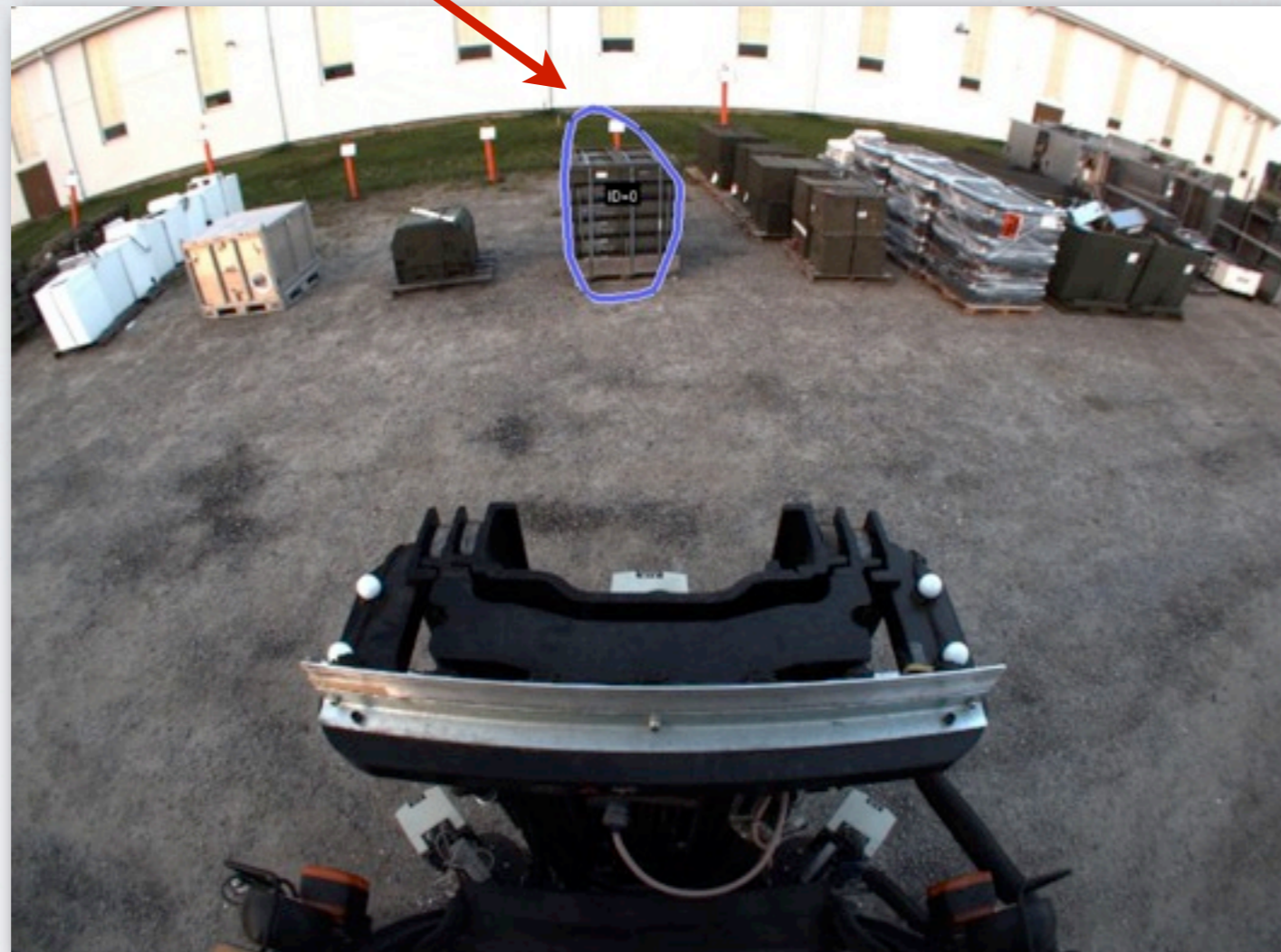
[ISER 2010; IJRR 2012]

Matthew Walter

Tuesday, February 5, 13

# Object Reacquisition

Matthew Walter

Tuesday, February 5, 13

# Object Reacquisition

Tuesday, February 5, 13

# Model Instantiation

User circles object in tablet image



Robot's forward-facing camera image

Matthew Walter

Tuesday, February 5, 13

# Model Instantiation



SIFT features extracted from initial image

Matthew Walter

# Model Instantiation



View 0 (user gesture)



SIFT features extracted from initial image

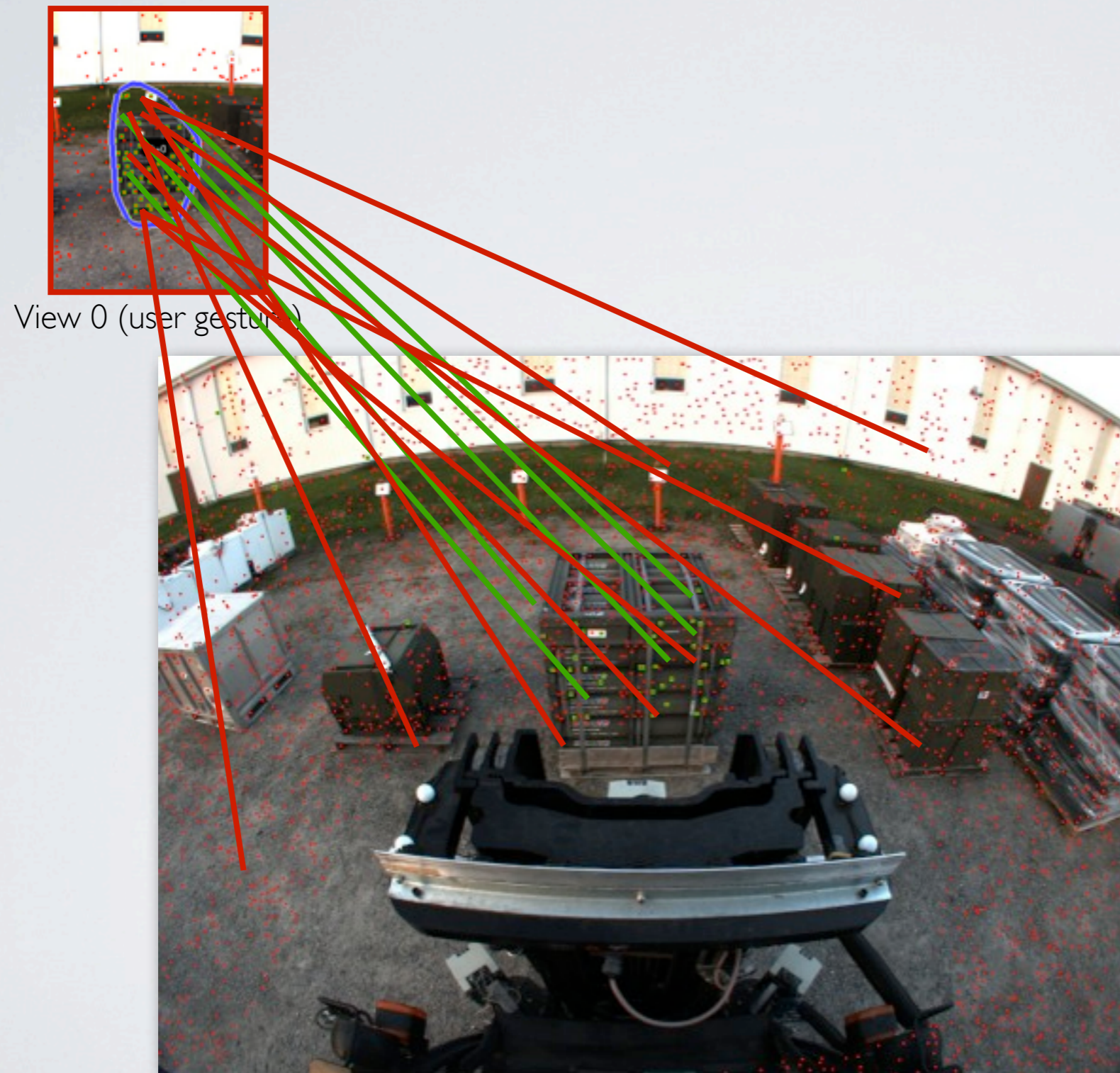Initialize model $\mathcal{M}_i$ to contain single view $v_{i1}$

Tuesday, February 5, 13

# Single-View Matching



View 0 (user gesture)



SIFT features extracted from new image

Extract features and match against all views

Matthew Walter
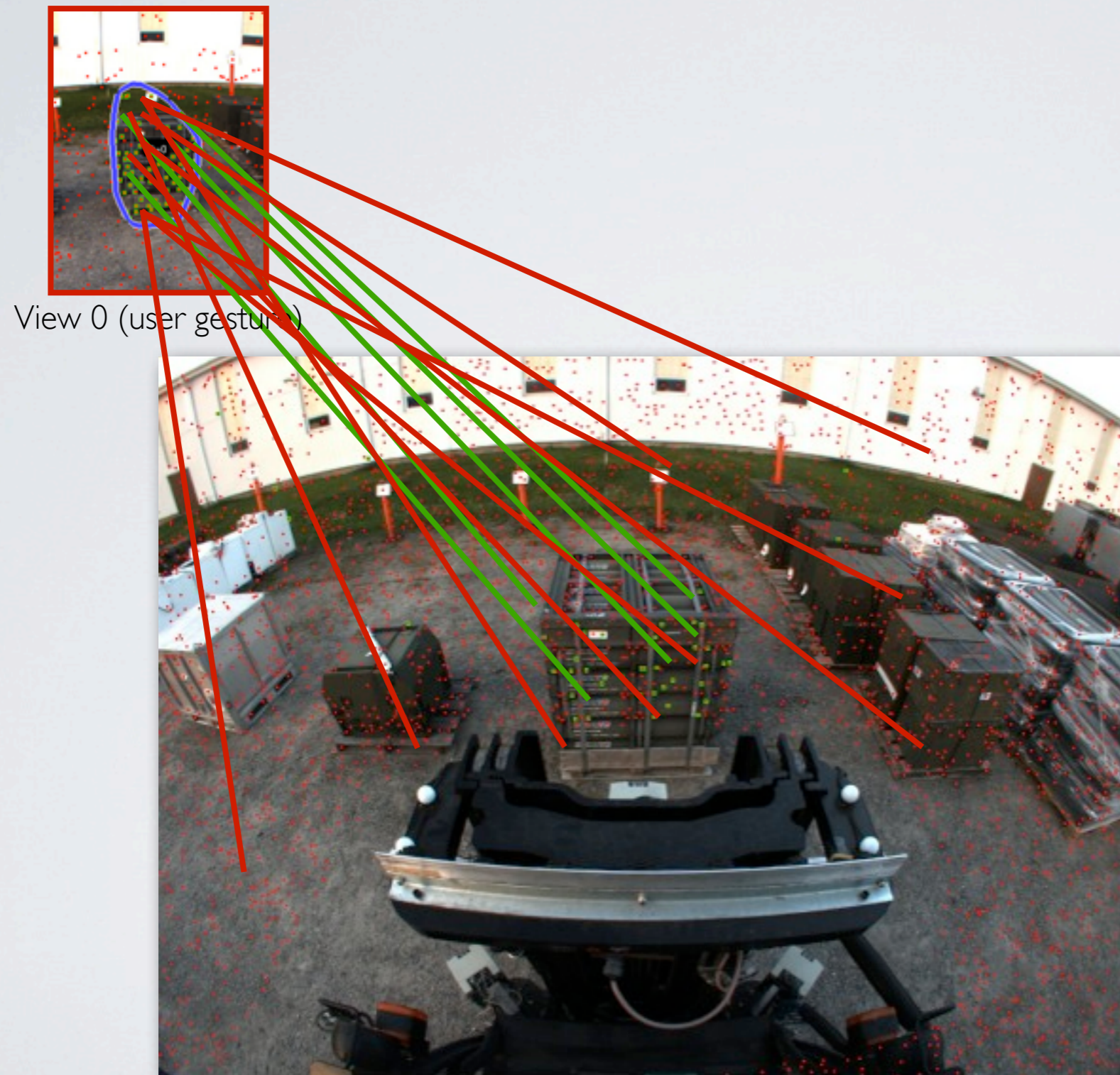
Tuesday, February 5, 13

# Single-View Matching



View 0 (user gesture)

SIFT features extracted from new image

Extract features and
match against
all views

Matthew Walter

Tuesday, February 5, 13

# Single-View Matching



View 0 (user gesture)

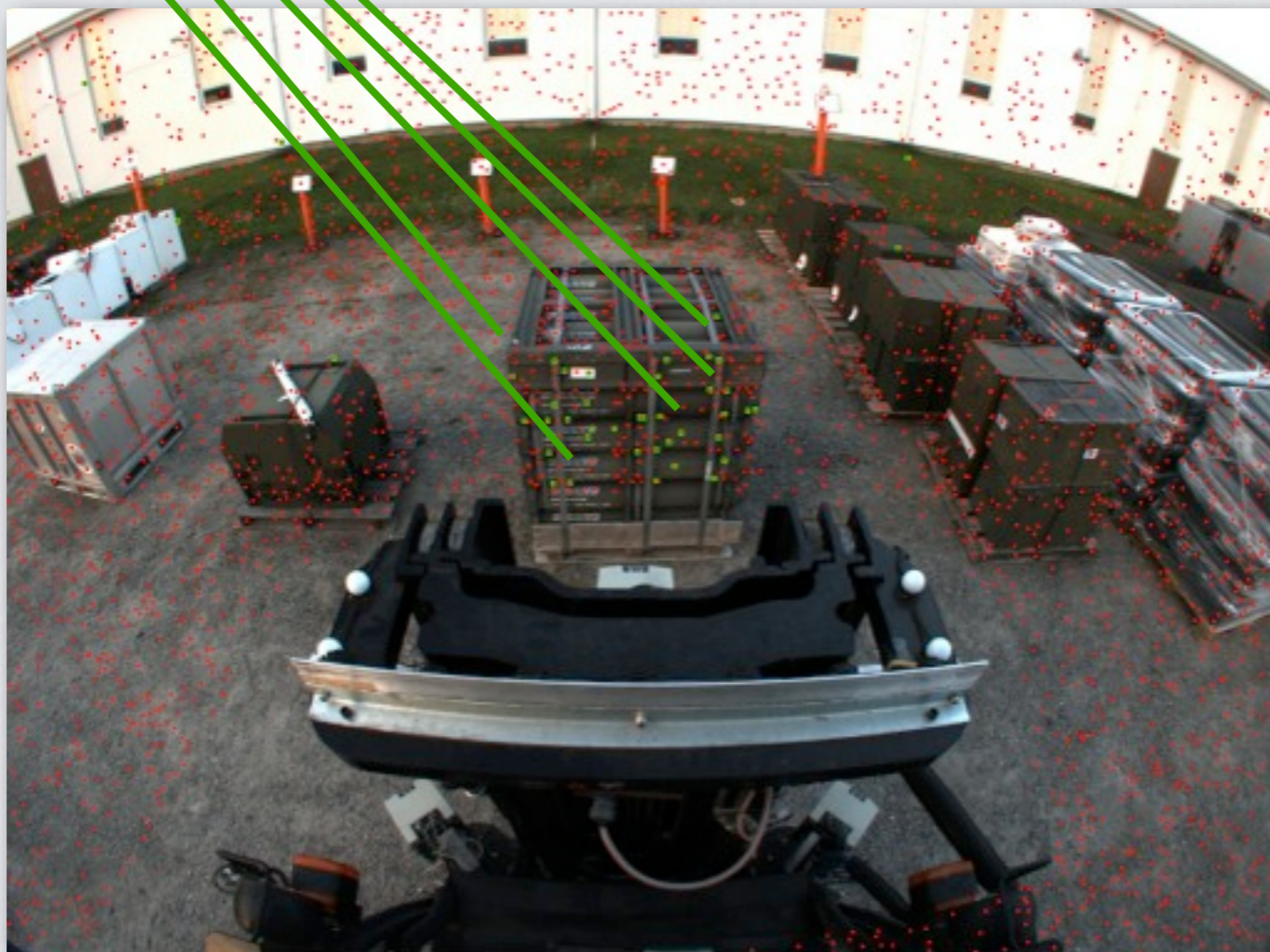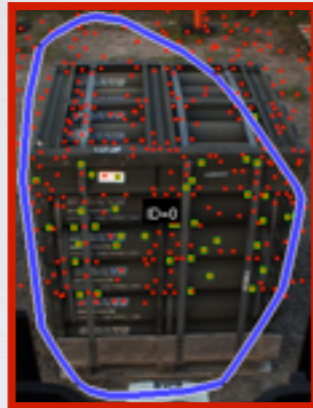SIFT features extracted from new image

## RANSAC

1. Sample a subset of pairs

2. Estimate corresponding image-to-image transformation (plane-projective homography)

3. Check consistency with other pairs

4. Repeat if inconsistent

Matthew Walter

Tuesday, February 5, 13

# Single-View Matching



View 0 (user gesture)

SIFT features extracted from new image

## RANSAC

1. Sample a subset of pairs

2. Estimate corresponding image-to-image transformation (plane-projective homography)

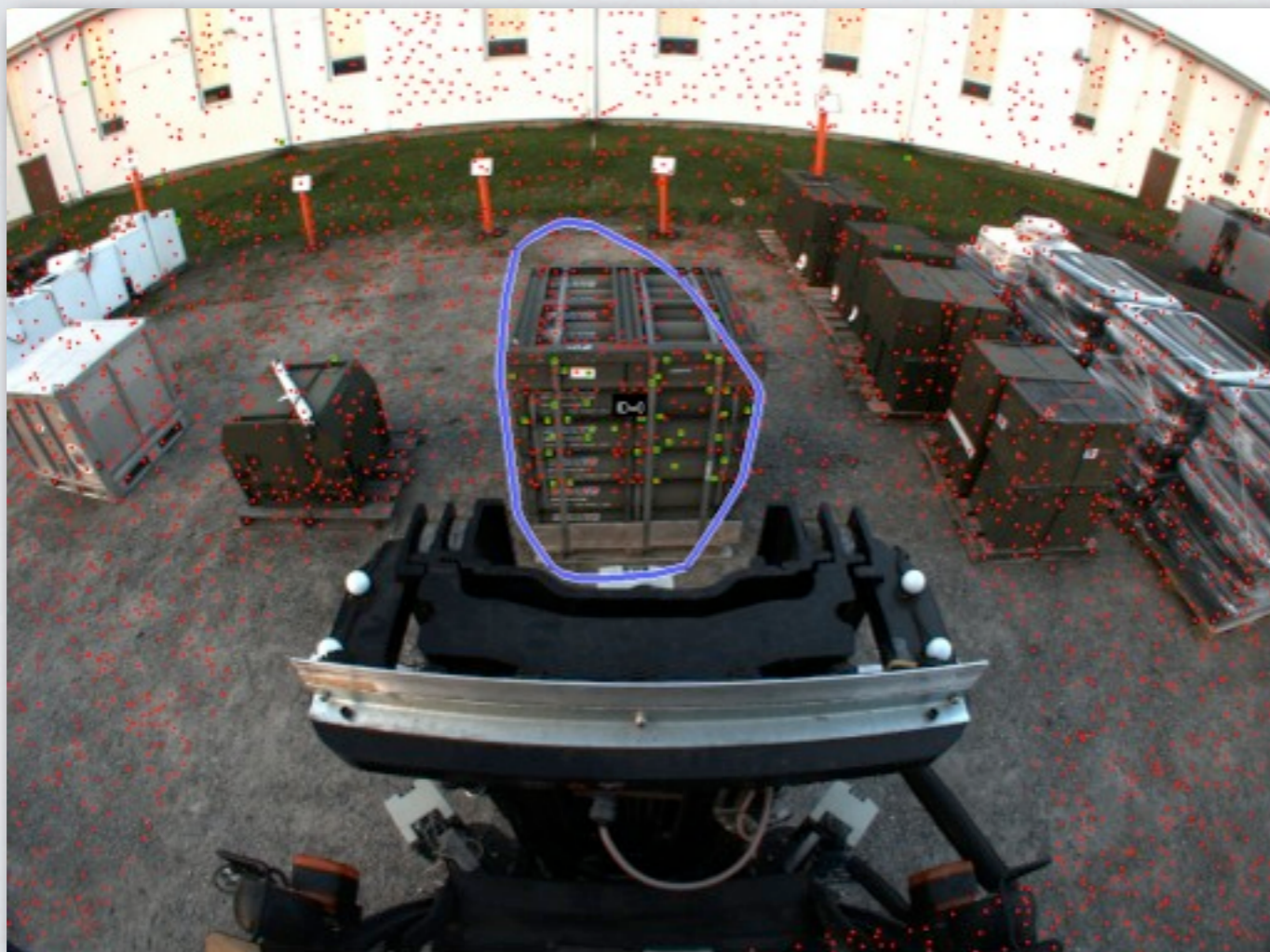3. Check consistency with other pairs

4. Repeat if inconsistent

Matthew Walter

Tuesday, February 5, 13

# Model Augmentation



View 0 (user gesture)          View 1

SIFT features extracted from new image

Generate segmentation and add new view

Matthew Walter

Tuesday, February 5, 13

# Model Augmentation



View 0 (user gesture)        View 1

SIFT features extracted from new image

Repeat as object appearance changes

Matthew Walter

Tuesday, February 5, 13

# Model Augmentation



View 0 (user gesture)      View 1      View 2



SIFT features extracted from new image

Repeat as object appearance changes

Matthew Walter

Tuesday, February 5, 13
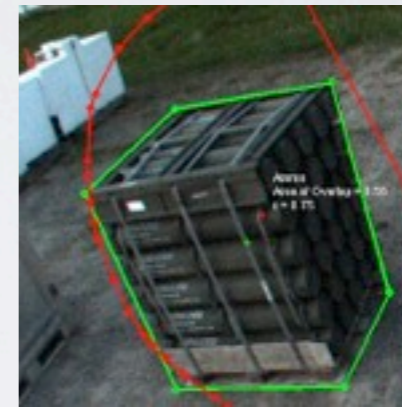
# One-shot Appearance Learning

Models opportunistically capture rich appearance variations



View 0 (user gesture)

Matthew Walter

Tuesday, February 5, 13

# One-shot Appearance Learning

## Models opportunistically capture rich appearance variations
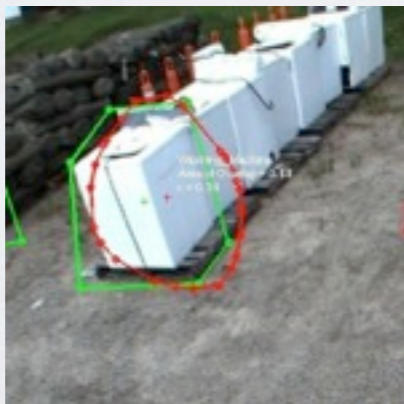


Model 1

| View 0 (user gesture) | View 1 | View 2 | View 3 | View 4 | View 5 |

Model 2

| View 0 (user gesture) | View 1 | View 2 | View 3 | View 4 | View 5 |

MIT   Matthew Walter

Tuesday, February 5, 13

# Visual Memory Results

- Active, outdoor military warehouse

- Tour and reacquisition separated by hours/days

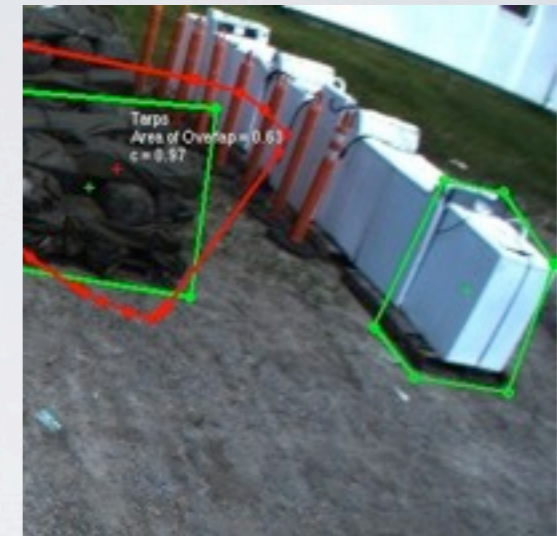- Training and detection with different cameras

- Varying conditions

Matthew Walter

Tuesday, February 5, 13

# Visual Memory Results



$$precision = \frac{TP}{TP + \underline{FP}}$$
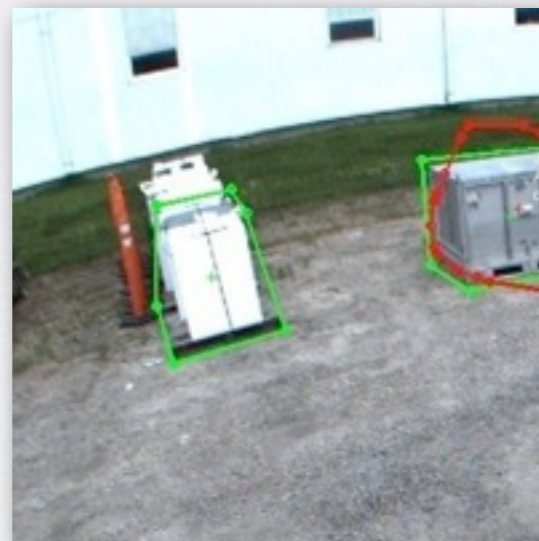
$$recall = \frac{TP}{TP + \underline{FN}}$$

| Scenario | Train | Test | Delta T | Precision | Recall |
|----------|-----------|---------|----------|-----------|--------|
| 1 | Afternoon | Afternoon | 5 min | 94% | 54% |
| 2 | Evening | Evening | 5 min | 100% | 95% |
| 3 | Morning | Evening | 14 hours | 100% | 93% |
| 4 | Morning | Evening | 10 hours | 100% | 94% |
| 5 | Noon | Evening | 7 hours | 100% | 94% |

Matthew Walter

Tuesday, February 5, 13

# Visual Memory Results

- Severe saturation

- Motion blur

- Unobserved viewpoints

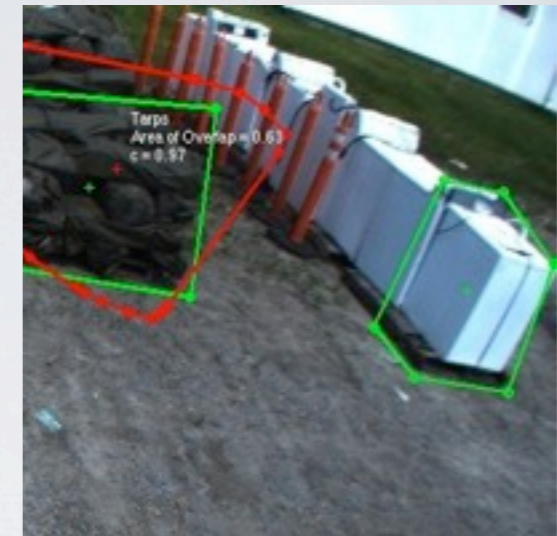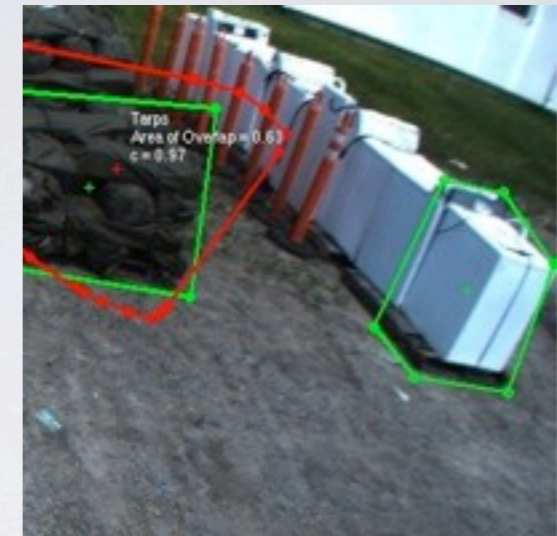

Training example

Saturation

Matthew Walter

Tuesday, February 5, 13

# Visual Memory Results



- Severe saturation

- Motion blur

- Unobserved viewpoints



Training example

Saturation

New viewpoint

Matthew Walter

Tuesday, February 5, 13
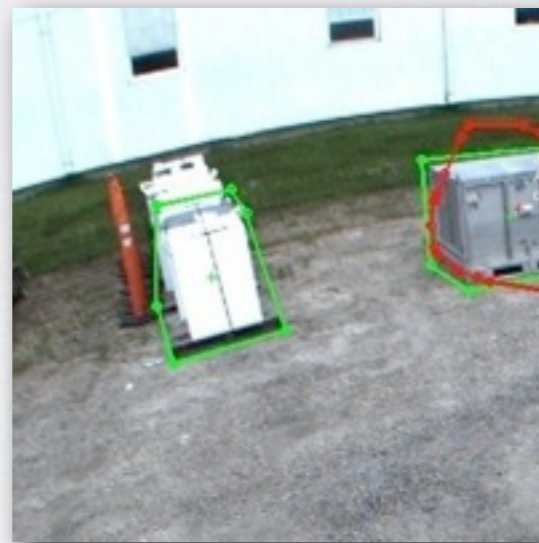
# Visual Memory Results



- Severe saturation

- Motion blur

- Unobserved viewpoints
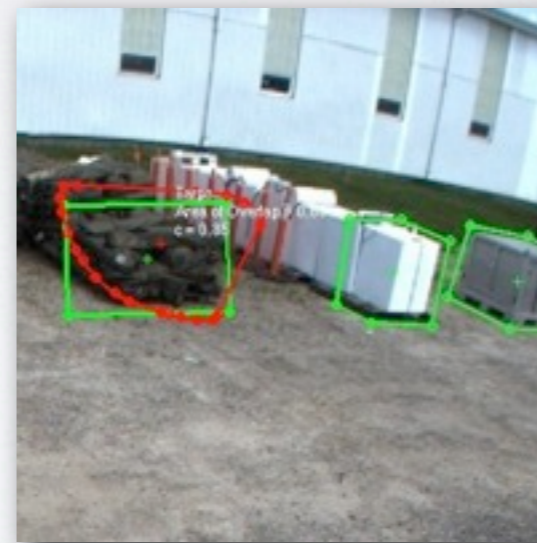
Training example | Saturation | New viewpoint | New viewpoint

Matthew Walter

Tuesday, February 5, 13

# Visual Memory Results

Matthew Walter

Tuesday, February 5, 13

# Visual Memory Results

Matthew Walter

Tuesday, February 5, 13

# Symbol Grounding Problem



Place the lifted tyre pallet, next to another tyre pallet on the trolley.

Lift the tire pallet in the air, then proceed to deposit it to the right of the tire pallet already on the table right in front of you.

Place the pallet of tires on the left side of the trailer.

Please lift the set of six tires up and set them on the trailer, to the right of the set of tires already on it.

lift the tire pallet you are carrying and set it on the truck in front of you

Place the pallet of tires that is on the forklift next to the pallet of tires that is already loaded on the trailer.

Lift tire pallet. Move to unoccupied location on truck. Lower tire pallet. Reverse to starting location. Lower forks. End.

Matthew Walter

# Symbol Grounding Problem

Linguistic elements $\longrightarrow$ Correct referents in the robot's world model

"Grounding"



Put the tire pallet on the truck

Place the lifted tyre pallet, next to another tyre pallet on the trolley.

Lift the tire pallet in the air, then proceed to deposit it to the right of the tire pallet already on the table right in front of you.

Place the pallet of tires on the left side of the trailer.

Please lift the set of six tires up and set them on the trailer, to the right of the set of tires already on it.

lift the tire pallet you are carrying and set it on the truck in front of you

Place the pallet of tires that is on the forklift next to the pallet of tires that is already loaded on the trailer.

Lift tire pallet. Move to unoccupied location on truck. Lower tire pallet. Reverse to starting location. Lower forks. End.

Matthew Walter

Tuesday, February 5, 13

# Symbol Grounding Problem

''Put the tire pallet on the truck''

- Objects
- Spatial relations
- Actions
- Places

- Object library
- Transformations, relative positions
- Paths, motion primitives, torques
- Positions, orientations

Matthew Walter

Tuesday, February 5, 13

# Symbol Grounding Problem

"Put the tire pallet on the truck"

- Objects
- Spatial relations
- Actions
- Places

- Object library
- Transformations, relative positions
- Paths, motion primitives, torques
- Positions, orientations

Matthew Walter

Tuesday, February 5, 13

# Symbol Grounding Problem

"Put the tire pallet on the truck"

- Objects

- Spatial relations

- Actions

- Places

- Object library

- Transformations, relative positions

- Paths, motion primitives, torques

- Positions, orientations

Matthew Walter

# Symbol Grounding Problem

"Put the tire pallet on the truck"

- Objects
- Spatial relations
- Actions
- Places

- Object library
- Transformations, relative positions
- Paths, motion primitives, torques
- Positions, orientations

Matthew Walter

# Symbol Grounding Problem

"Put the tire pallet on the truck"

- Objects
- Spatial relations
- Actions
- Places



- Object library
- Transformations, relative positions
- Paths, motion primitives, torques
- Positions, orientations

Matthew Walter

Tuesday, February 5, 13

# Symbol Grounding Problem

## "Put the tire pallet on the truck"

- Objects
- Spatial relations
- Actions
- Places

- Object library
- Transformations, relative positions
- Paths, motion primitives, torques
- Positions, orientations

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

$$\underset{\text{groundings}}{\arg\max}\ p\,(\text{groundings}|\text{language})$$

objects, actions, relations, places                "Drive to the tire pallet"
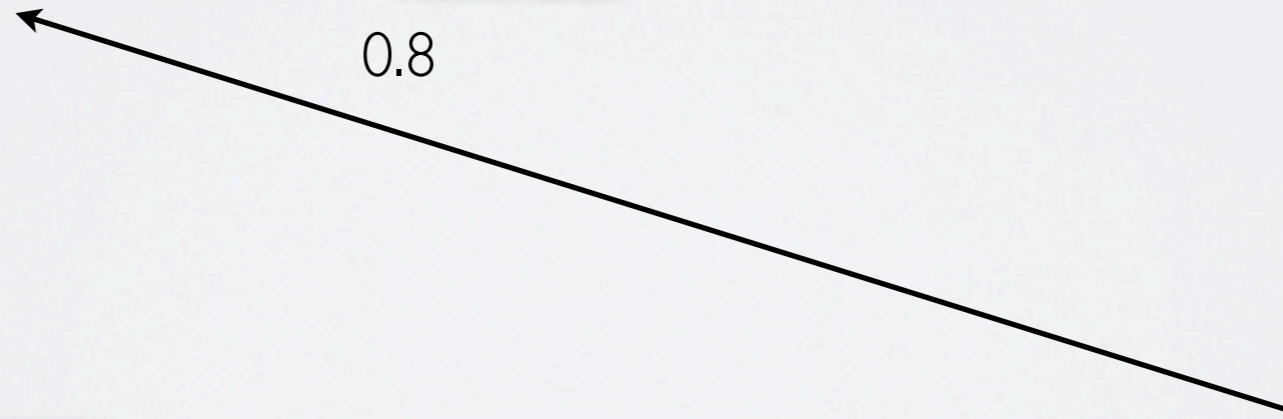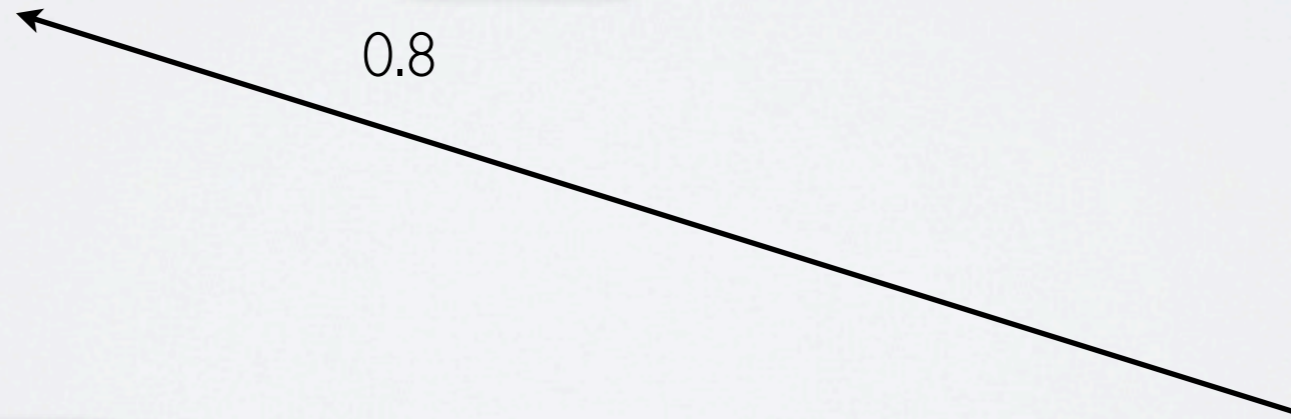
Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg\max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg\max_{\gamma \in \mathcal{X}} p\left(\gamma|\lambda\right)$$

0.1

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg\max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$



0.8

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg\max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$



0.1

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg \max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$



0.8

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

$$\underset{\text{groundings}}{\arg\max}\; p\,(\text{groundings}|\text{language})$$

objects, actions, relations, places     "Put the tire pallet on the truck"

[AAAI 2011; AI Magazine 2011]

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

"Put the tire pallet on the truck"

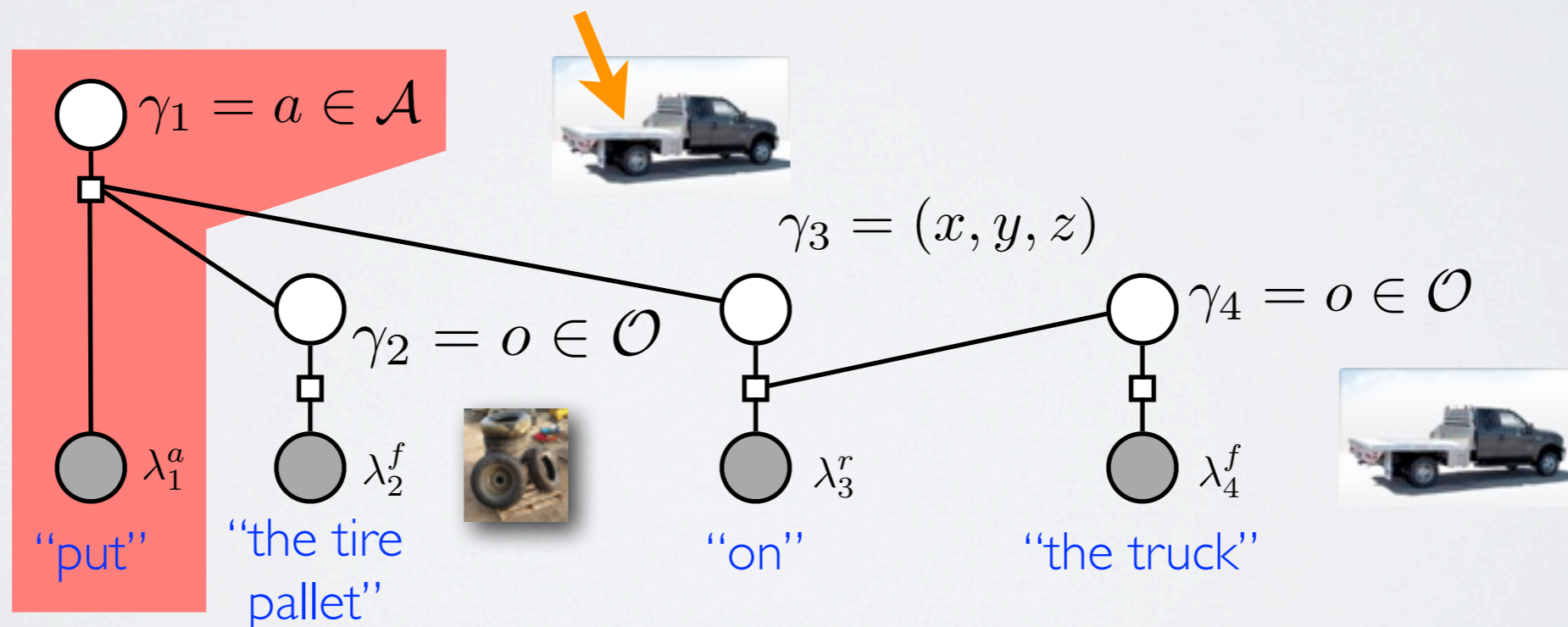$$\arg\max_{\Gamma} \ (\gamma_1, \gamma_2, \gamma_3, \gamma_4 | \lambda)$$



[AAAI 2011; AI Magazine 2011]

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

"Put **the tire pallet** on the truck"

$$\arg \max_{\Gamma} \left( \gamma_1, \gamma_2, \gamma_3, \gamma_4 | \lambda \right)$$



$\gamma_1 = a \in \mathcal{A}$

$\gamma_3 = (x, y, z)$

$\gamma_2 = o \in \mathcal{O}$

$\gamma_4 = o \in \mathcal{O}$

$\lambda_1^a$    $\lambda_2^f$    $\lambda_3^r$    $\lambda_4^f$

"put"    "the tire pallet"    "on"    "the truck"

MIT    Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

"Put the tire pallet on the truck"

$$\arg\max_{\Gamma} \; (\gamma_1, \gamma_2, \gamma_3, \gamma_4 | \lambda)$$



$\gamma_1 = a \in \mathcal{A}$

$\gamma_3 = (x, y, z)$

$\gamma_2 = o \in \mathcal{O}$

$\gamma_4 = o \in \mathcal{O}$

$\lambda_1^a$     $\lambda_2^f$     $\lambda_3^r$     $\lambda_4^f$

"put"    "the tire pallet"    "on"    "the truck"

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

"Put the tire pallet **on** the truck"

$$\arg\max_{\Gamma} \; (\gamma_1, \gamma_2, \gamma_3, \gamma_4 | \lambda)$$



$\gamma_1 = a \in \mathcal{A}$

$\gamma_3 = (x, y, z)$

$\gamma_4 = o \in \mathcal{O}$

$\gamma_2 = o \in \mathcal{O}$

$\lambda_1^a$

$\lambda_2^f$

$\lambda_3^r$

$\lambda_4^f$

"put"      "the tire pallet"      "on"      "the truck"

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

"**Put** the tire pallet on the truck"

$$\arg\max_{\Gamma} \ (\gamma_1, \gamma_2, \gamma_3, \gamma_4 | \lambda)$$



$\gamma_1 = a \in \mathcal{A}$

$\gamma_3 = (x, y, z)$

$\gamma_2 = o \in \mathcal{O}$

$\gamma_4 = o \in \mathcal{O}$

$\lambda_1^a$    $\lambda_2^f$    $\lambda_3^r$    $\lambda_4^f$

"put"    "the tire pallet"    "on"    "the truck"

Matthew Walter

Tuesday, February 5, 13

I. Importance of Situational Awareness

II. Persistent Object Awareness with Vision

**III. Semantic Map Learning from Natural Language Descriptions**

IV. Future Directions

V. Conclusions

Matthew Walter

# Beyond Objects to Spaces

- Going beyond metric maps

- Human-centric representations of space
  - Spatial relations
  - Semantic attributes (names, use, etc.)
  - Connectivity

Matthew Walter

Tuesday, February 5, 13

# State-of-the-Art in Semantic Mapping

- Spatial Semantic Hierarchy (Kuipers 2000)

- Augment SLAM metric/topological SLAM maps with semantic layers



Courtesy: Zender et al. 2008

- Infer semantic properties from multiple modalities:

  - Object recognition (Zender et al. 2008; Pronobis et al. 2020)

  - Spoken descriptions and other supervised labels
    (Diosi et al. 2005; Zender et al. 2008; Pronobis et al. 2020)

  - Place classification (Zender et al. 2008; Pronobis et al. 2020)

Matthew Walter

Tuesday, February 5, 13

# Building Semantic Maps with Natural Language

- Learn knowledge representation from narrated tour

- Challenges:
  - People convey high-level concepts but robot perception is low-level
  - Spoken descriptions are ambiguous

Matthew Walter

# Building Semantic Maps with Natural Language

- Solution:
  - Joint metric, topologic, & semantic model supports information fusion
  - Efficient inference strategy
  - Enable layers to influence one another

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t)$$

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs



$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t)$$

Topology $G_t = (V_t, E_t)$

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs



$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t)$$

Topology $G_t = (V_t, E_t)$

Vertex poses

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs



$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t)$$

Topology $G_t = (V_t, E_t)$

Semantic labels

Vertex poses

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs



Sensor stream

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t)$$

Topology $G_t = (V_t, E_t)$

Semantic labels

Vertex poses

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs

Sensor stream    Odometry

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t)$$

Topology $G_t = (V_t, E_t)$    Semantic labels

Vertex poses

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs
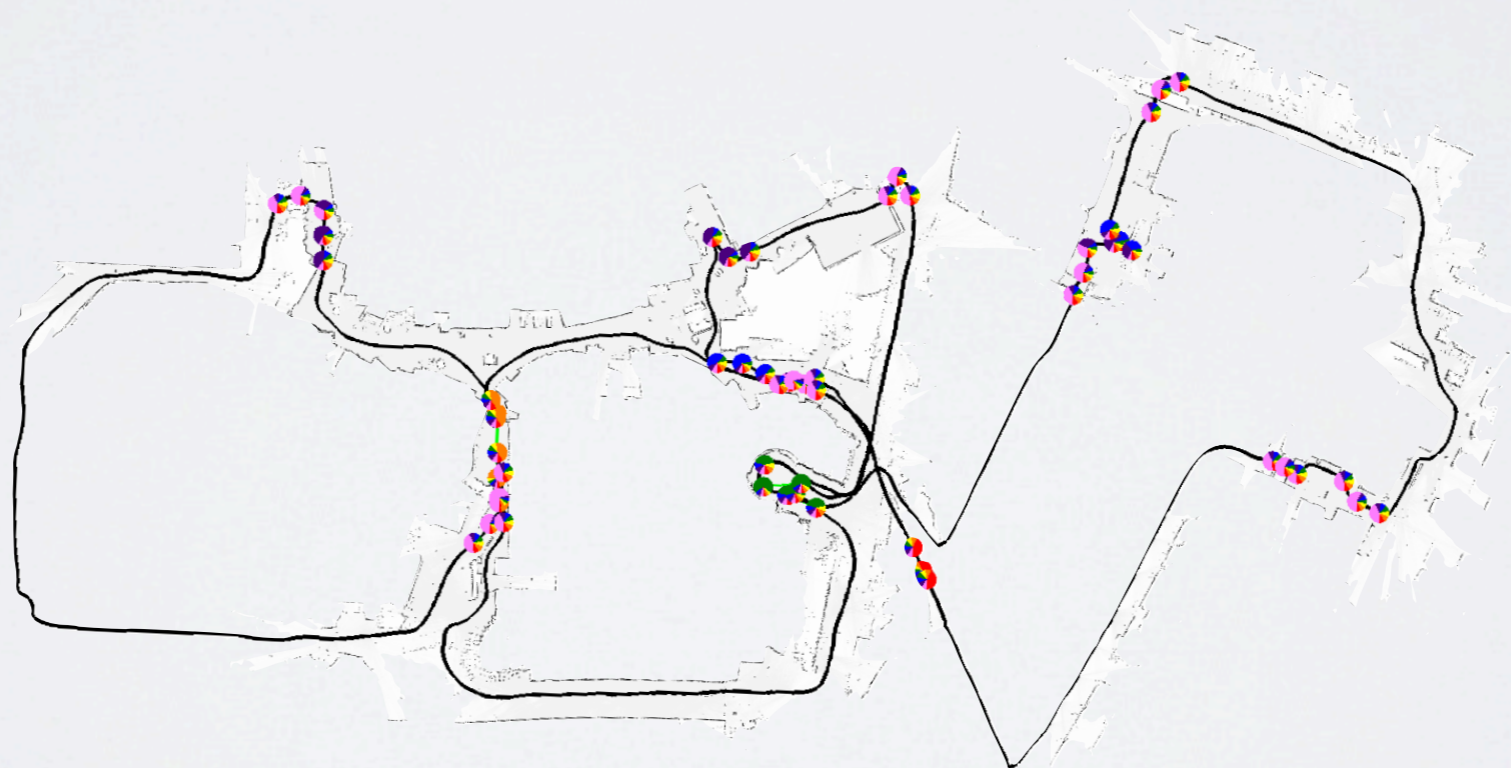
Matthew Walter

# Model: Posterior over Semantic Graphs

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t) = p(L_t | X_t, G_t, z^t, u^t, \lambda^t) \ p(X_t | G_t, z^t, u^t, \lambda^t) \ p(G_t | z^t, u^t, \lambda^t)$$

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t) = p(L_t | X_t, G_t, z^t, u^t, \lambda^t) \; p(X_t | G_t, z^t, u^t, \lambda^t) \boxed{p(G_t | z^t, u^t, \lambda^t)}$$

Sample-based
representation

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t) = p(L_t | X_t, G_t, z^t, u^t, \lambda^t) \boxed{p(X_t | G_t, z^t, u^t, \lambda^t)} p(G_t | z^t, u^t, \lambda^t)$$

$$\boxed{\begin{array}{c|c} \text{Gaussian} & \text{Sample-based} \\ \text{(information form)} & \text{representation} \end{array}}$$

$$p(X_t | G_t, z^t, u^t, \lambda^t) = \mathcal{N}^{-1}(X_t; \Sigma_t^{-1}, \eta_t)$$

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t) = \boxed{p(L_t | X_t, G_t, z^t, u^t, \lambda^t)} \; p(X_t | G_t, z^t, u^t, \lambda^t) \; p(G_t | z^t, u^t, \lambda^t)$$

Dirichlet

Gaussian
(information form)

Sample-based
representation

Matthew Walter

Tuesday, February 5, 13

# Model: Posterior over Semantic Graphs

$$p(G_t, X_t, L_t | z^t, u^t, \lambda^t) = p(L_t | X_t, G_t, z^t, u^t, \lambda^t) \ p(X_t | G_t, z^t, u^t, \lambda^t) \ p(G_t | z^t, u^t, \lambda^t)$$

Dirichlet

Gaussian
(information form)

Sample-based
representation

Matthew Walter

Tuesday, February 5, 13

# Rao-Blackwellized Particle Filter

Input: $P_{t-1} = \left\{ G_{t-1}^{(i)}, X_{t-1}^{(i)}, L_{t-1}^{(i)} w_{t-1}^{(i)} \right\}$   $(u_t, z_t, \lambda_t)$
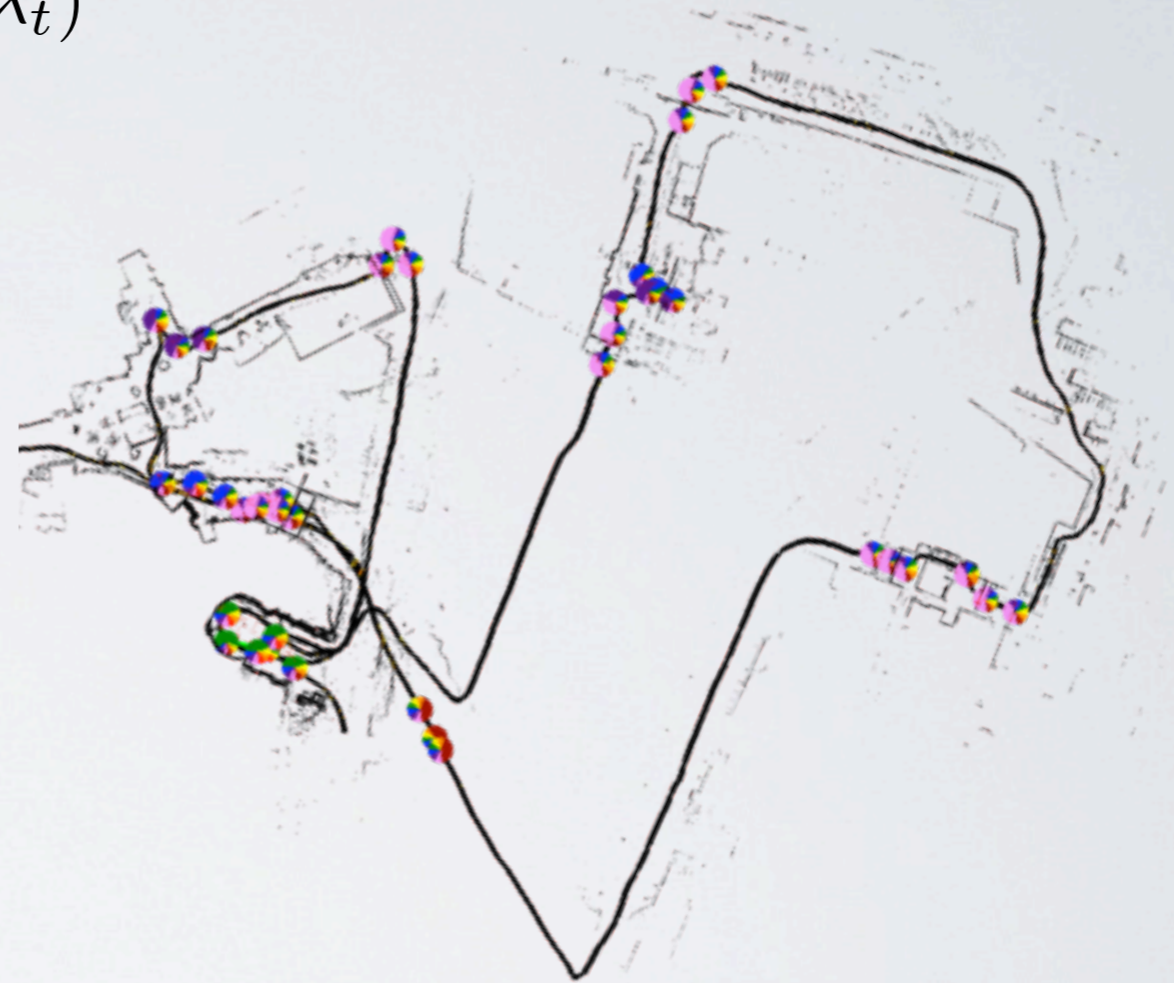
for each particle i

1  Propose modifications to topology based on metric and semantic maps

2  Perform Bayesian update of Gaussian
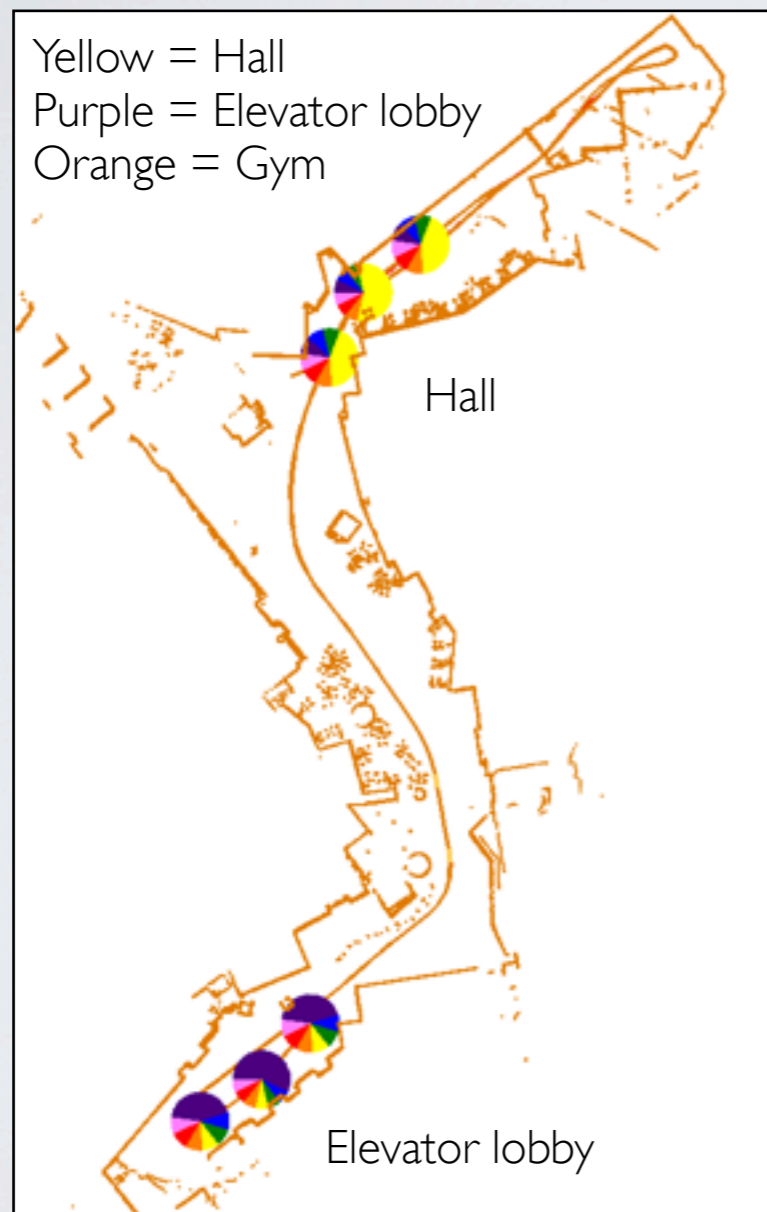
3  Update Dirichlet over labels based on language

4  Update weights based on metric observations

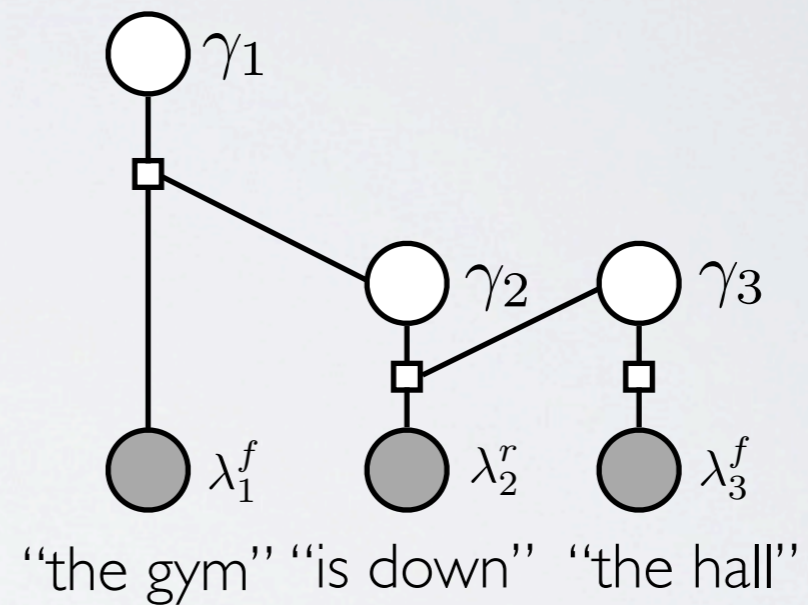Return: $P_t^{(i)} = \left\{ G_t^{(i)}, X_t^{(i)}, L_t^{(i)} w_t^{(i)} \right\}$

Matthew Walter

Tuesday, February 5, 13

# Rao-Blackwellized Particle Filter

Input: $P_{t-1} = \left\{ G_{t-1}^{(i)}, X_{t-1}^{(i)}, L_{t-1}^{(i)} w_{t-1}^{(i)} \right\} \quad (u_t, z_t, \lambda_t)$

for each particle i

1  Propose modifications to topology based on metric and semantic maps

2  Perform Bayesian update of Gaussian

3  Update Dirichlet over labels based on language

4  Update weights based on metric observations

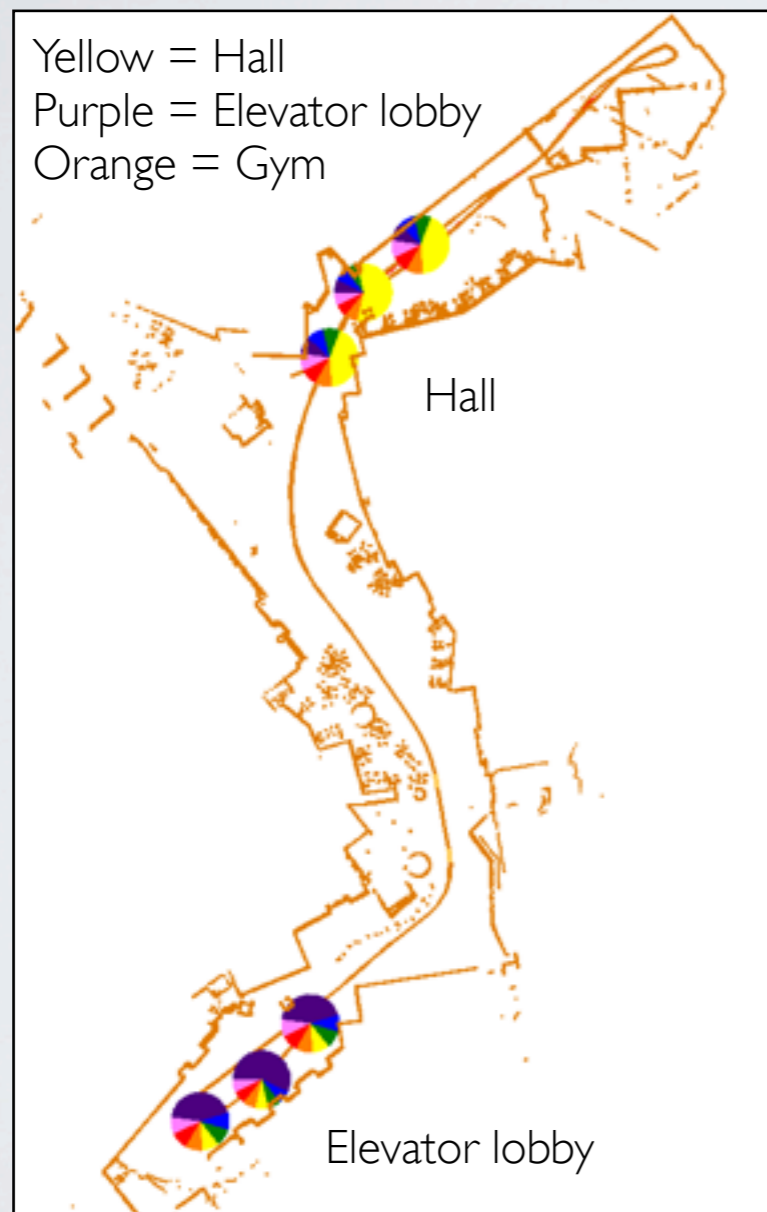Return: $P_t^{(i)} = \left\{ G_t^{(i)}, X_t^{(i)}, L_t^{(i)} w_t^{(i)} \right\}$

Matthew Walter

Tuesday, February 5, 13

# Rao-Blackwellized Particle Filter

Input: $P_{t-1} = \left\{ G_{t-1}^{(i)}, X_{t-1}^{(i)}, L_{t-1}^{(i)} w_{t-1}^{(i)} \right\}$  $(u_t, z_t, \lambda_t)$

for each particle i

1  Propose modifications to topology based on metric and semantic maps

2  Perform Bayesian update of Gaussian

3  Update Dirichlet over labels based on language

4  Update weights based on metric observations

Return: $P_t^{(i)} = \left\{ G_t^{(i)}, X_t^{(i)}, L_t^{(i)} w_t^{(i)} \right\}$

Matthew Walter

Tuesday, February 5, 13

# Rao-Blackwellized Particle Filter

Input: $P_{t-1} = \left\{ G_{t-1}^{(i)}, X_{t-1}^{(i)}, L_{t-1}^{(i)} w_{t-1}^{(i)} \right\}$   $(u_t, z_t, \lambda_t)$

for each particle i

1  Propose modifications to topology based on metric and semantic maps

2  Perform Bayesian update of Gaussian

3  Update Dirichlet over labels based on language

4  Update weights based on metric observations

Return: $P_t^{(i)} = \left\{ G_t^{(i)}, X_t^{(i)}, L_t^{(i)} w_t^{(i)} \right\}$

"the gym is down the hall"

Yellow = Hall
Purple = Elevator lobby
Orange = Gym

Matthew Walter

Tuesday, February 5, 13

# Incorporating Natural Language Descriptions

"the gym is down the hall"



Yellow = Hall
Purple = Elevator lobby
Orange = Gym

Hall

Elevator lobby

Matthew Walter

Tuesday, February 5, 13

# Incorporating Natural Language Descriptions

"the gym is down the hall"

Yellow = Hall
Purple = Elevator lobby
Orange = Gym

Hall

Elevator lobby

$\gamma_1$

$\gamma_2$  $\gamma_3$

$\lambda_1^f$   $\lambda_2^r$   $\lambda_3^f$

"the gym" "is down" "the hall"

$$p(L_t^{(i)}|L_{t-1}^{(i)}, G_t^{(i)}, X_t^{(i)}, \lambda_t) =$$
$$\sum_\gamma p(L_t^{(i)}|\gamma, L_{t-1}^{(i)}, \lambda_t) \times p(\gamma|L_{t-1}^{(i)}, G_t^{(i)}, X_t^{(i)}, \lambda_t)$$

Matthew Walter

# Incorporating Natural Language Descriptions

Matthew Walter

Tuesday, February 5, 13

# No Language Constraints

Tuesday, February 5, 13

# No Language Constraints

MIT          Matthew Walter

Tuesday, February 5, 13

# No Language Constraints

Matthew Walter

Tuesday, February 5, 13

# No Language Constraints

Tuesday, February 5, 13

# With Language Constraints

MIT

Matthew Walter

Tuesday, February 5, 13

# No Language Constraints

MIT

Tuesday, February 5, 13

# No Language Constraints



Legend:
- Gym
- Elevator lobby
- Courtyard
- Cafeteria
- Hallway
- Amphiteater
- Entrance

20 m

Matthew Walter

# With Language Constraints



Legend:
- Gym
- Elevator lobby
- Courtyard
- Cafeteria
- Hallway
- Amphiteater
- Entrance

Loop closures

Loop closure

20 m

Matthew Walter

Tuesday, February 5, 13

# Preliminary Results - With Language Constraints



Guide: Good afternoon, Please follow me
Robot: Following

Matthew Walter

Tuesday, February 5, 13

Matthew Walter

# Enhancing Models of Objects and Space

- Object category recognition
  - Data-driven models
  - Transfer learning
  - Limited supervision via human intervention
  - Efficient retrieval and matching

- New sources of information
  - Objects (e.g., co-occurrence)
  - Vision-based scene classification
  - Higher-level concepts
  - Building topology databases

- Exploration-based natural language grounding

Matthew Walter

Tuesday, February 5, 13

# Where are We Going?

People

Robots

- Objects

- Places

- Actions

- People

- Events



Images



Laser scans



Wheel torques

Figure 2 The unimate PUMA 562 robot arm



Joint angles

Matthew Walter

Tuesday, February 5, 13

# Where are We Going?

People

Robots

- Objects

- Places

- Actions

- People

- Events


Images


Laser scans


Wheel torques


Joint angles
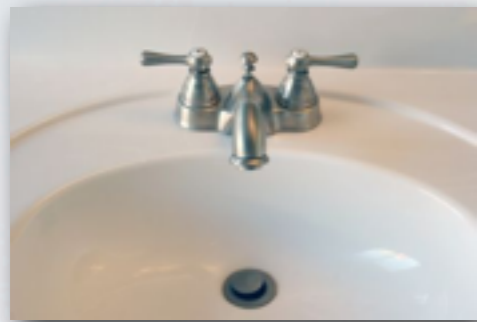
Matthew Walter

Tuesday, February 5, 13

# Learning Rich Action Spaces

- Low-level, object-specific actions don't scale

- Long-term planning & inference is intractable
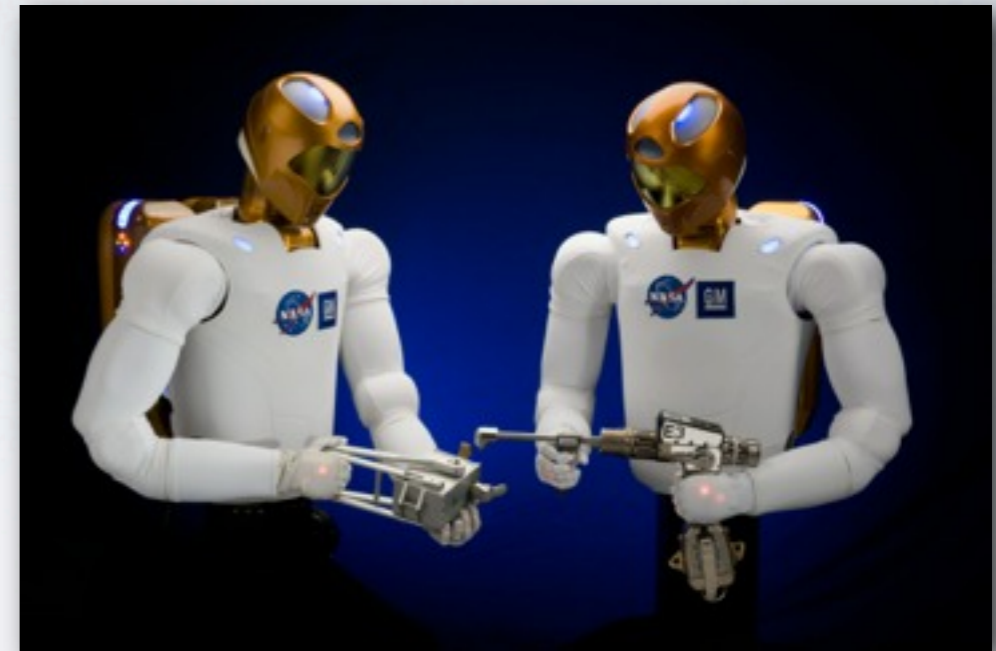




Courtesy: Willow Garage

Tuesday, February 5, 13

# Learning Rich Action Spaces

- Low-level, object-specific actions don't scale

- Long-term planning & inference is intractable

Matthew Walter

Tuesday, February 5, 13

# Learning Rich Action Spaces

- Robots need higher-level representations
  - Structured state/action space
  - Affordance-based action model
  - Affordances are grounded in perception

- Human-provided information is critical to efficient learning

- Robots must formulate representation based on experience

Matthew Walter
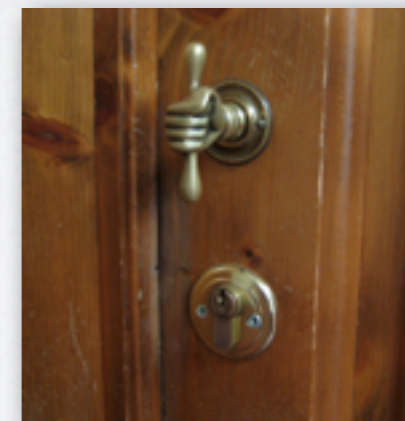
Tuesday, February 5, 13

# Learning Rich Action Spaces

- Robots need higher-level representations
  - Structured state/action space
  - Affordance-based action model
  - Affordances are grounded in perception

- Human-provided information is critical to efficient learning

- Robots must formulate representation based on experience

Matthew Walter

Tuesday, February 5, 13

# Learning Rich Action Spaces

- Robots need higher-level representations
  - Structured state/action space
  - Affordance-based action model
  - Affordances are grounded in perception

- Human-provided information is critical to efficient learning

- Robots must formulate representation based on experience

Matthew Walter

# Learning via Deliberate Actions

- Robots need higher-level representations
  - Structured state/action space
  - Affordance-based action model
  - Affordances are grounded in perception

- Human-provided information is critical to efficient learning

- Robots must formulate representation based on experience

Matthew Walter

Tuesday, February 5, 13

Matthew Walter

Tuesday, February 5, 13

# Contributions

- Showed that human-robot collaboration requires intuitive control

- Argued that the key missing capability is situational awareness

- Perception is critical to enabling awareness

- Demonstrated algorithms that opportunistically learn rich models of objects and space from human-provided cues
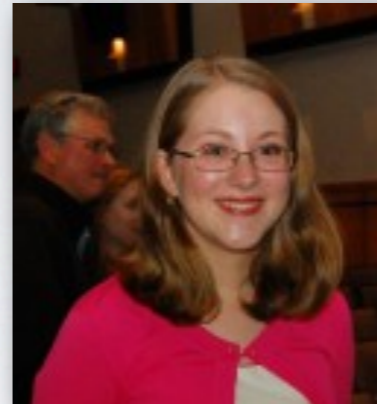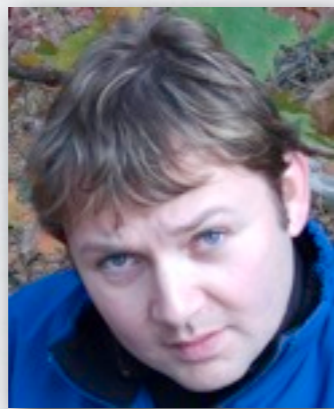
Matthew Walter

Tuesday, February 5, 13

# Contributions


Sachi Hemachandra


Stefanie Tellex


Bianca Homberg


Sudeep Pillai


Yuli Friedman


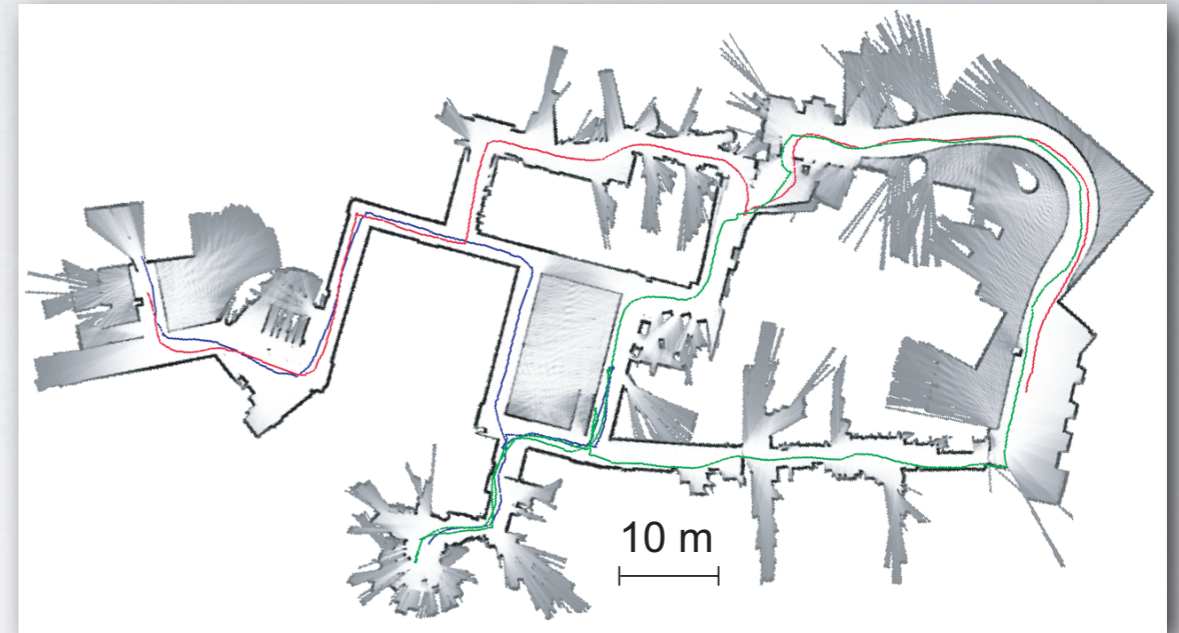Matthew Antone


Seth Teller

Matthew Walter

Tuesday, February 5, 13

# Contributions

- Showed that human-robot collaboration requires intuitive control

- Argued that the key missing capability is situational awareness

- Perception is critical to enabling awareness

- Demonstrated algorithms that opportunistically learn rich models of objects and space from human-provided cues

Matthew Walter

Tuesday, February 5, 13

# Getting There

- Autonomous navigation

- Motion planning & control

- Planning under uncertainty

- Manipulation

- Localization & mapping

- Perception

- Efficient control

- Natural interaction

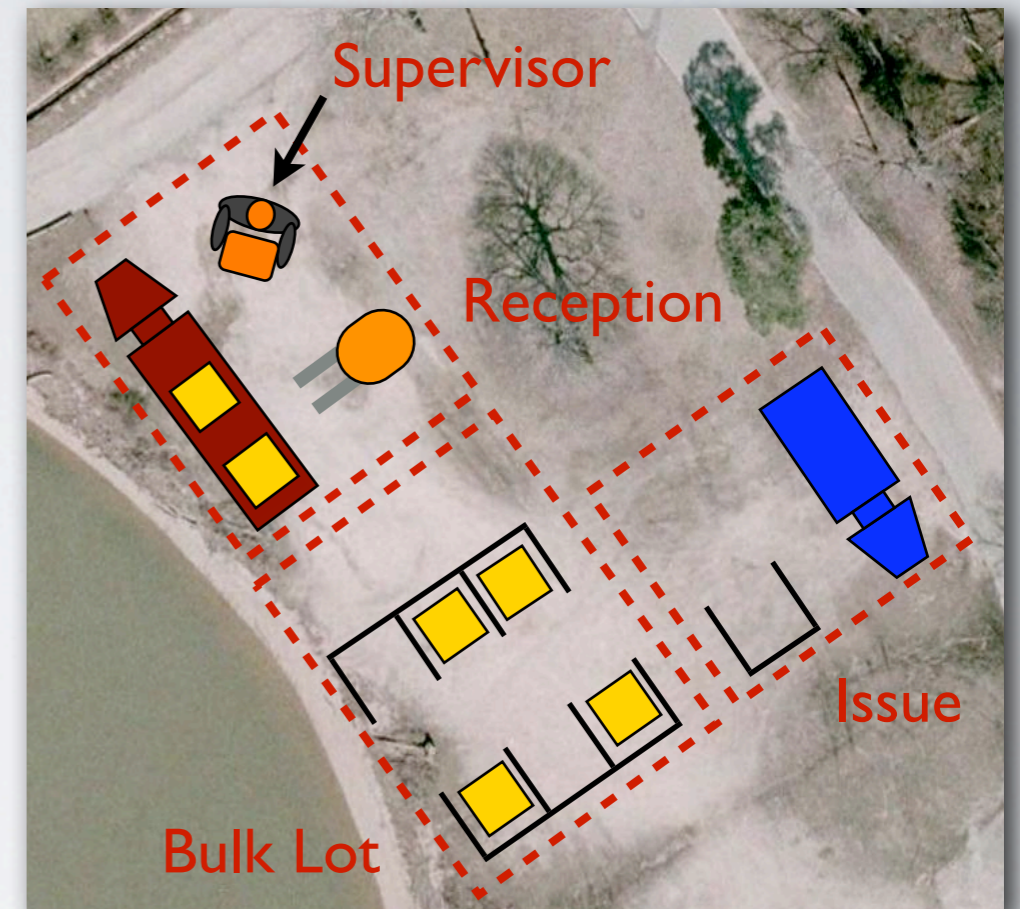- Trusted autonomy



Navigation using uncalibrated cameras



RRT



Anytime RRT*: optimal planning

Matthew Walter

Tuesday, February 5, 13

# Mobile Manipulation for Logistics



1. Pick up objects off trucks or the ground

2. Transport items to storage locations

3. Load particular objects onto customer trucks or the ground

Matthew Walter
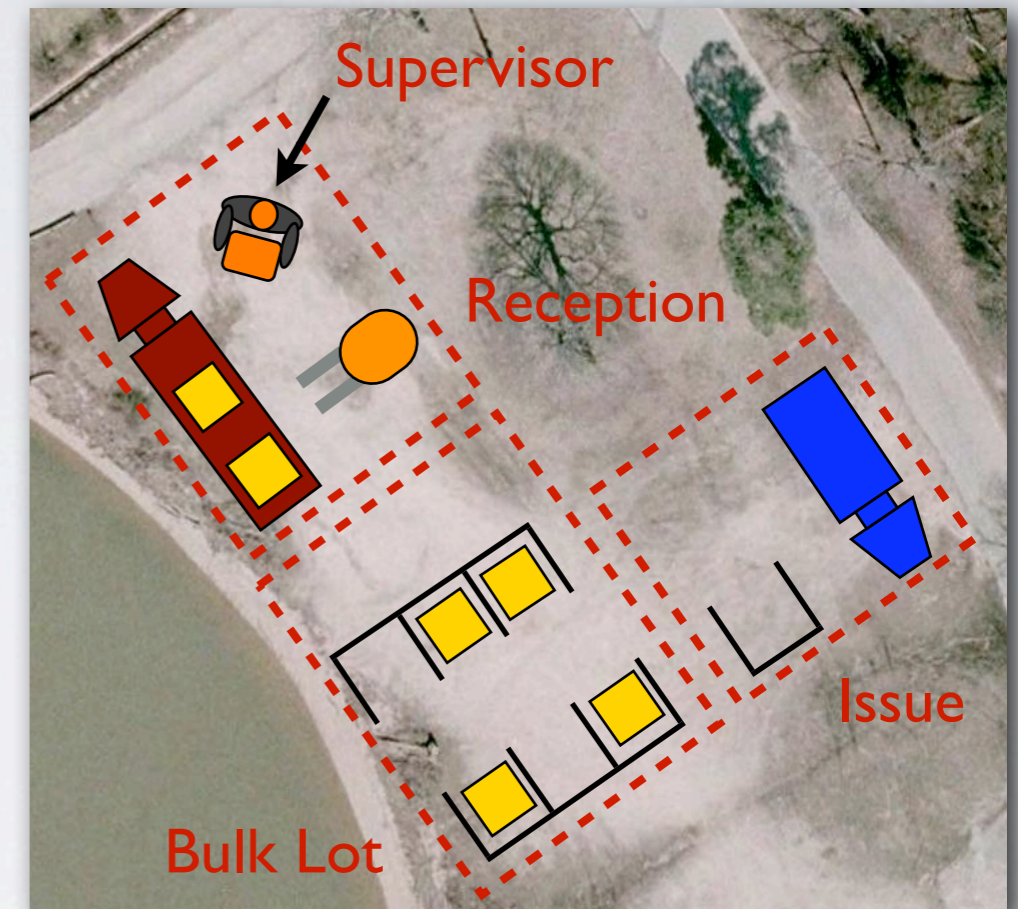
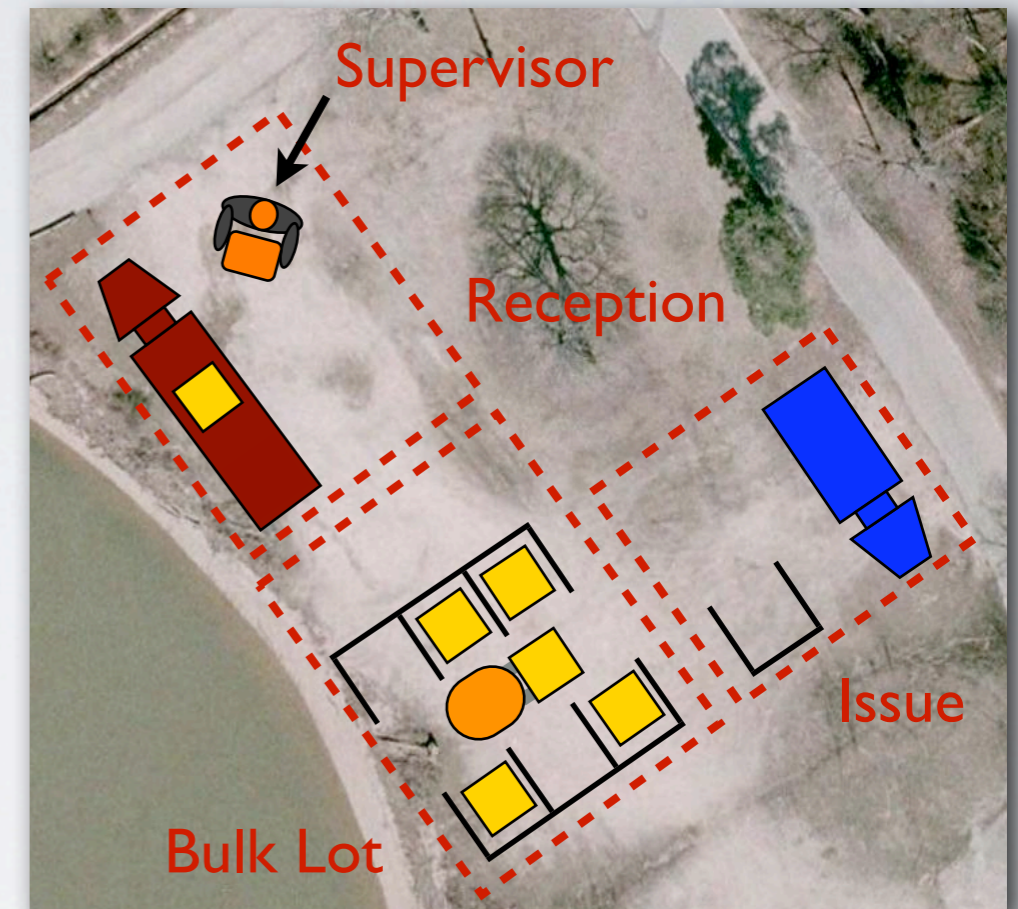Tuesday, February 5, 13

# Mobile Manipulation for Logistics





1. Pick up objects off trucks or the ground

2. Transport items to storage locations

3. Load particular objects onto customer trucks or the ground

Matthew Walter

# Mobile Manipulation for Logistics



1. Pick up objects off trucks or the ground

2. Transport items to storage locations

3. Load particular objects onto customer trucks or the ground

Matthew Walter

# Mobile Manipulation for Logistics
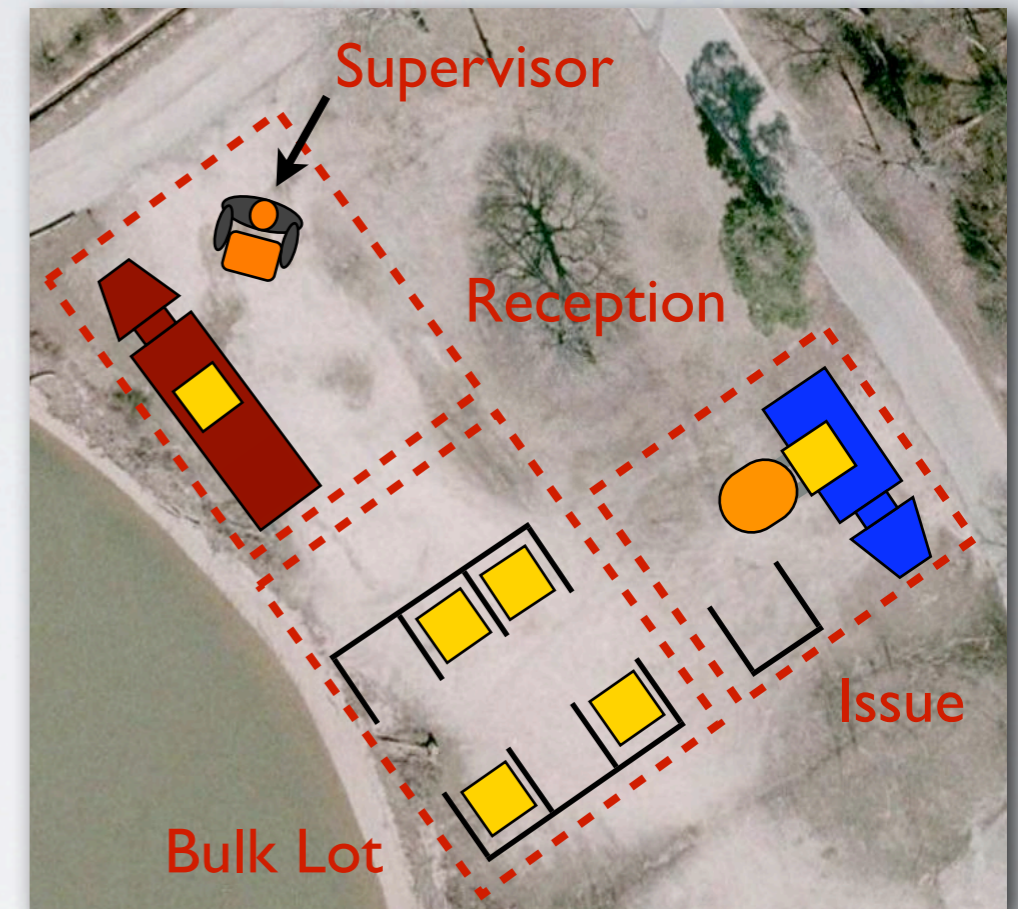


1. Pick up objects off trucks or the ground

2. Transport items to storage locations

3. Load particular objects onto customer trucks or the ground

Matthew Walter

# The Platform



Microphones          Cameras



Laser range-finders

Matthew Walter

# Teleoperation?

Matthew Walter

Tuesday, February 5, 13

# Teleoperation?

Matthew Walter

Tuesday, February 5, 13

# Shared Autonomy



- Hierarchical task-level autonomy
  - Reduce tasks into simpler sub-tasks

- Shared situational awareness

- Robot can request help when needed

Tuesday, February 5, 13

# Command via Shared World Model



- Hand-held tablet interface
  - Robot's eye view with annotated images
  - Onboard speech recognition
  - Interprets pen-based gestures

- Microphones pick up external speech

Matthew Walter

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

$$\underset{\text{groundings}}{\arg\max}\; p\,(\text{groundings}|\text{language})$$

objects, actions, relations, places       "Drive to the tire pallet"

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg \max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$

Matthew Walter
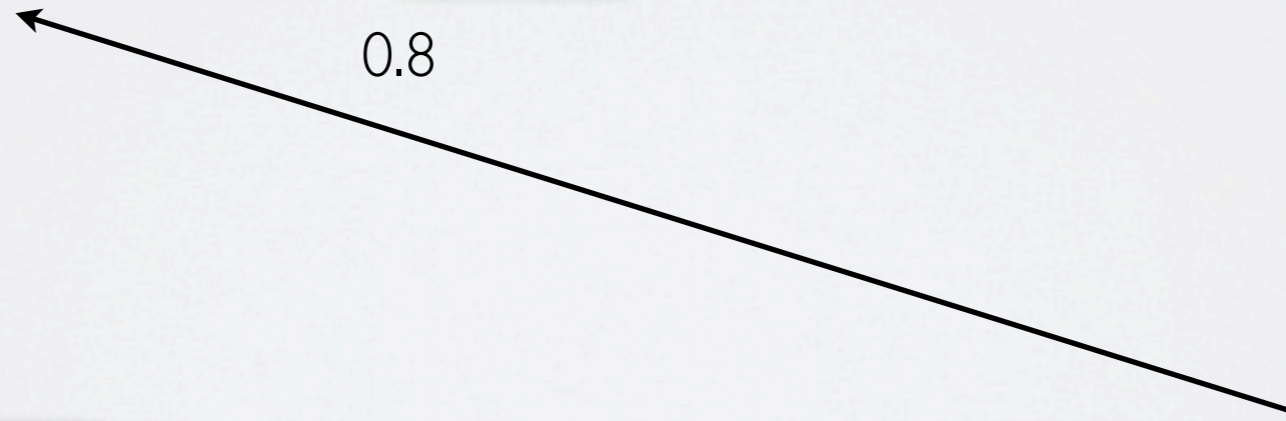
Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg\max_{\gamma\in\mathcal{X}} p\left(\gamma|\lambda\right)$$



0.1

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg \max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$



0.8

Matthew Walter
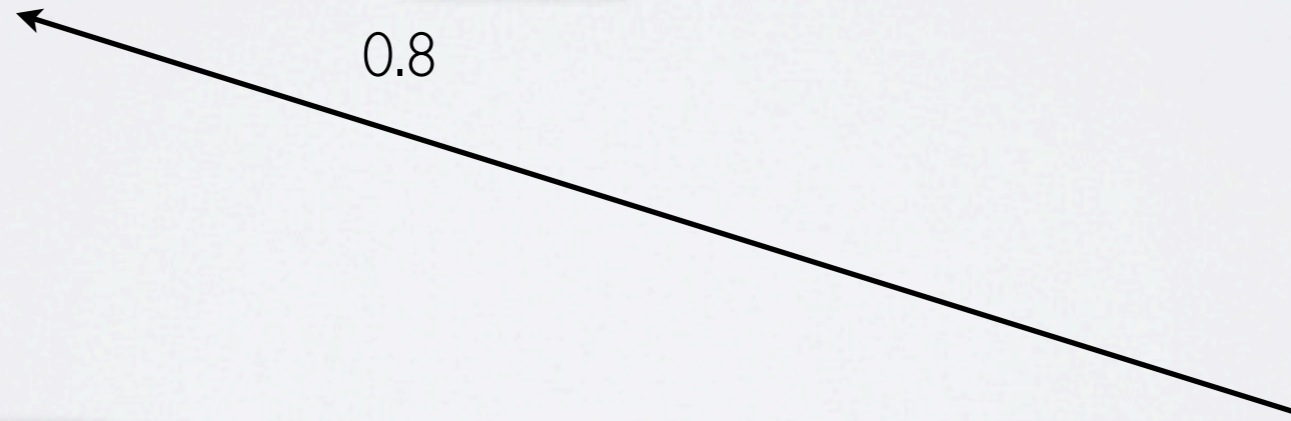
Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg \max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$



0.1

Matthew Walter

Tuesday, February 5, 13

# Grounding Natural Language Speech

(collaboration with S. Tellex, T. Kollar, S. Teller, & N. Roy)

"To the tire pallet"

$$\arg\max_{\gamma \in \mathcal{X}} p\left(\gamma | \lambda\right)$$



0.8

Matthew Walter

Tuesday, February 5, 13

# Learning the Grounding Distributions



Training Set

"To the tire pallet"

$\gamma = (x_1, y_1, z_1)$

!("To the tire pallet")

$\gamma = (x_2, y_2, z_2)$
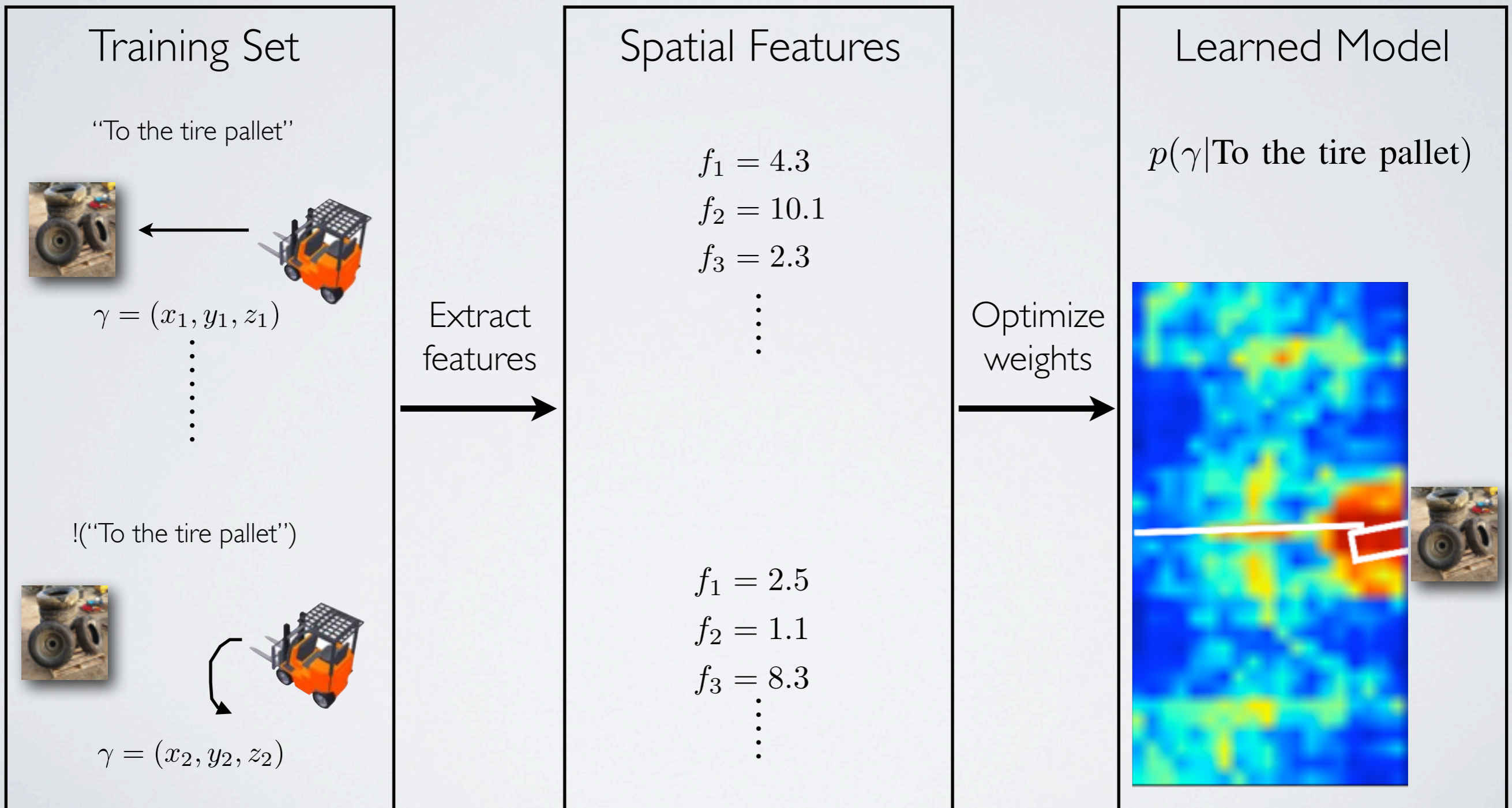
Matthew Walter

Tuesday, February 5, 13

# Learning the Grounding Distributions

# Learning the Grounding Distributions

Matthew Walter

Tuesday, February 5, 13

# Learning the Grounding Distributions

Tuesday, February 5, 13