

Semantic Perception under Uncertainty with Viewpoint-Dependent Models

Yuri Feldman

Under the supervision of Assoc. Prof. **Vadim Indelman**

Ph.D. Seminar, March 2022

Work partially supported by



Intro – Robot Autonomy

Key components:

Perception (Situational Awareness, Data Fusion)

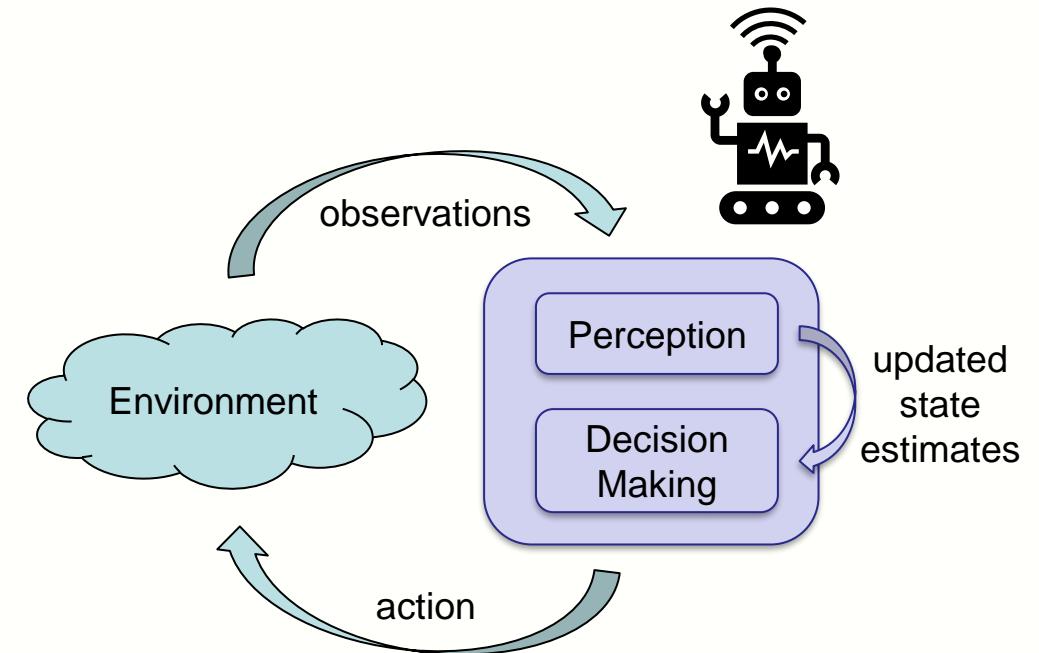
- Understanding of the environment and robot state within it

Decision Making

- Plan actions (towards task of interest)

Need to deal with **uncertainty**

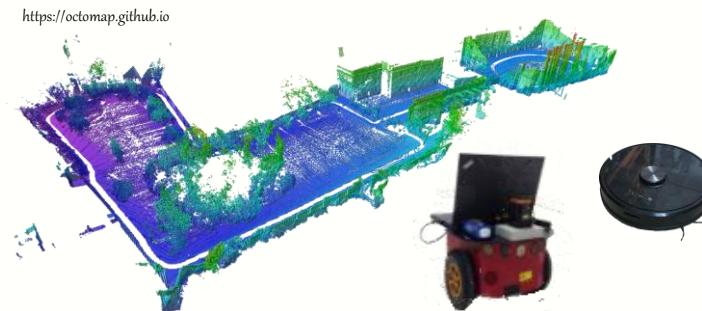
- Due to: noisy and aliased measurements, partial information, *imperfect models*...



Intro 2 – Semantic Perception

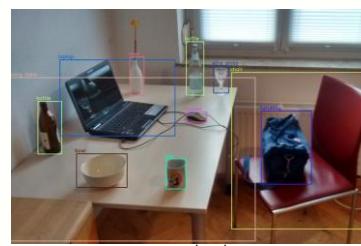
Geometric perception

- established methods exist
- SLAM – Simultaneous Localization and Mapping



Semantic perception

- required for less-structured tasks
- need resilience to per-frame errors for safe and reliable operation.



Outline

Viewpoint-Dependent models for Semantic Perception under Uncertainty

1. The semantic perception problem (Object-Level SLAM)
2. Viewpoint-dependent semantic measurement models
3. Contributions:
 - I. Classification under Model and Localization Uncertainty
 - II. Data Association-Aware Semantic Mapping and Localization
 - III. Semantic Perception with a Continuous Learned Representation

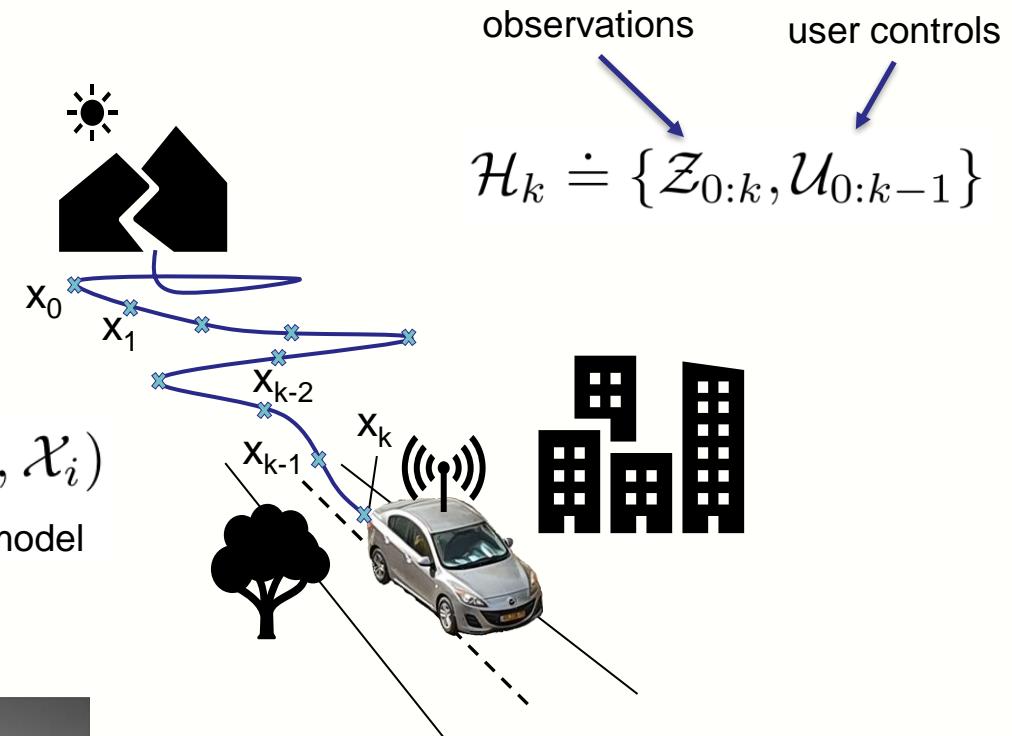
Intro 3 – (Geometric) Simultaneous Localization and Mapping (SLAM)

$$\begin{aligned} \mathcal{X}_{0:k}^*, \mathcal{L}^* &= \arg \max_{\mathcal{X}, \mathcal{L}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{L} \mid \mathcal{H}_k) \\ &= \arg \max_{\mathcal{X}, \mathcal{L}} \eta \cdot \mathbb{P}(\mathcal{X}_0) \prod_{\text{prior}} \mathbb{P}(\mathcal{X}_{i+1} \mid \mathcal{U}_i, \mathcal{X}_i) \prod_i \mathbb{P}(\mathcal{X}_i \mid \mathcal{X}_{0:i}) \prod_{j \in \mathcal{M}_i} \mathbb{P}(\mathcal{Z}_j \mid \mathcal{L}_j, \mathcal{X}_i) \end{aligned}$$

robot track landmarks

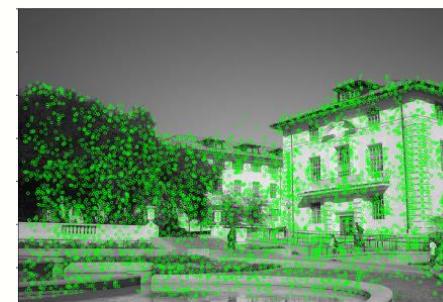
$$= \arg \max_{\mathcal{X}, \mathcal{L}} \eta \cdot \mathbb{P}(\mathcal{X}_0) \prod_{\text{prior}} \mathbb{P}(\mathcal{X}_{i+1} \mid \mathcal{U}_i, \mathcal{X}_i) \prod_i \mathbb{P}(\mathcal{X}_i \mid \mathcal{X}_{0:i}) \prod_{j \in \mathcal{M}_i} \mathbb{P}(\mathcal{Z}_j \mid \mathcal{L}_j, \mathcal{X}_i)$$

observations user controls



See for example:

Kaess et al. 2008 TRO (iSAM)
 Kaess et al. 2012 IJRR (iSAM2)
 And related publications



Intro 3 – (Geometric) Simultaneous Localization and Mapping (SLAM)

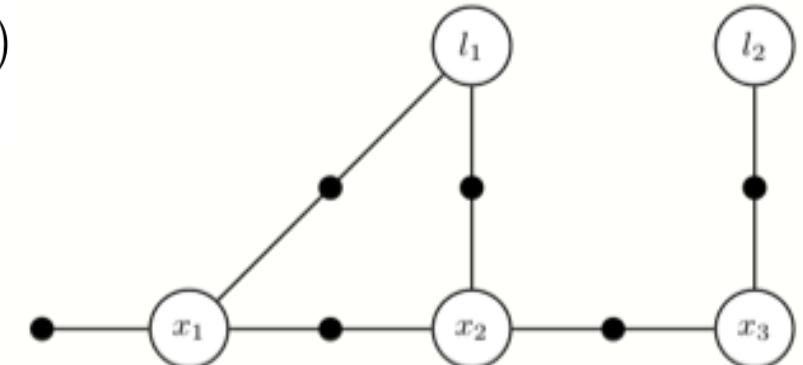
$$\mathcal{X}_{0:k}^*, \mathcal{L}^* = \arg \max_{\mathcal{X}, \mathcal{L}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{L} \mid \mathcal{H}_k)$$

robot track landmarks

$$= \arg \max_{\mathcal{X}, \mathcal{L}} \eta \cdot \mathbb{P}(\mathcal{X}_0) \prod_{\text{prior}} \prod_i \mathbb{P}(\mathcal{X}_{i+1} \mid \mathcal{U}_i, \mathcal{X}_i) \prod_{j \in \mathcal{M}_i} \mathbb{P}(\mathcal{Z}_j \mid \mathcal{L}_j, \mathcal{X}_i)$$

motion model observation model

Factor graph:



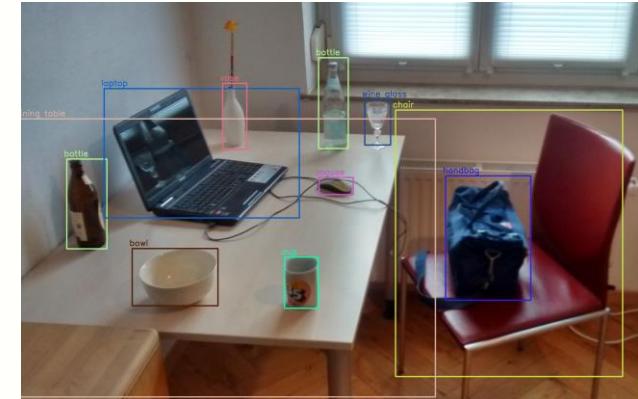
© gtsam tutorial

Semantic Perception: Object-Level SLAM

(Salas-Moreno et al. 13' CVPR, Choudhary et al. 14' IROS, Bowman et al. 17' ICRA, McCormack et al. 18 3dv, Nicholson 19' ral, Yang 19' TRO, ...)

$$\arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{C}, \mathcal{O} \mid \mathcal{H}_k)$$

object categories (discrete!) measurement and control history
robot track object geometry



© Wikipedia

$$= \arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{O} \mid \mathcal{C}, \mathcal{H}_k) \cdot \mathbb{P}(\mathcal{C} \mid \mathcal{H}_k)$$

(continuous) hypothesis hypothesis weight

$$c \in \{1, \dots, N\}$$

$$|\{(c_1, \dots, c_m)\}| = N^m$$

Semantic Perception: Object-Level SLAM

(Salas-Moreno et al. 13' CVPR, Choudhary et al. 14' IROS, Bowman et al. 17' ICRA, McCormack et al. 18 3dv, Nicholson 19' ral, Yang 19' TRO, ...)

$$\arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{C}, \mathcal{O} \mid \mathcal{H}_k)$$

object categories
(discrete!)

robot track

measurement
and control history

object geometry

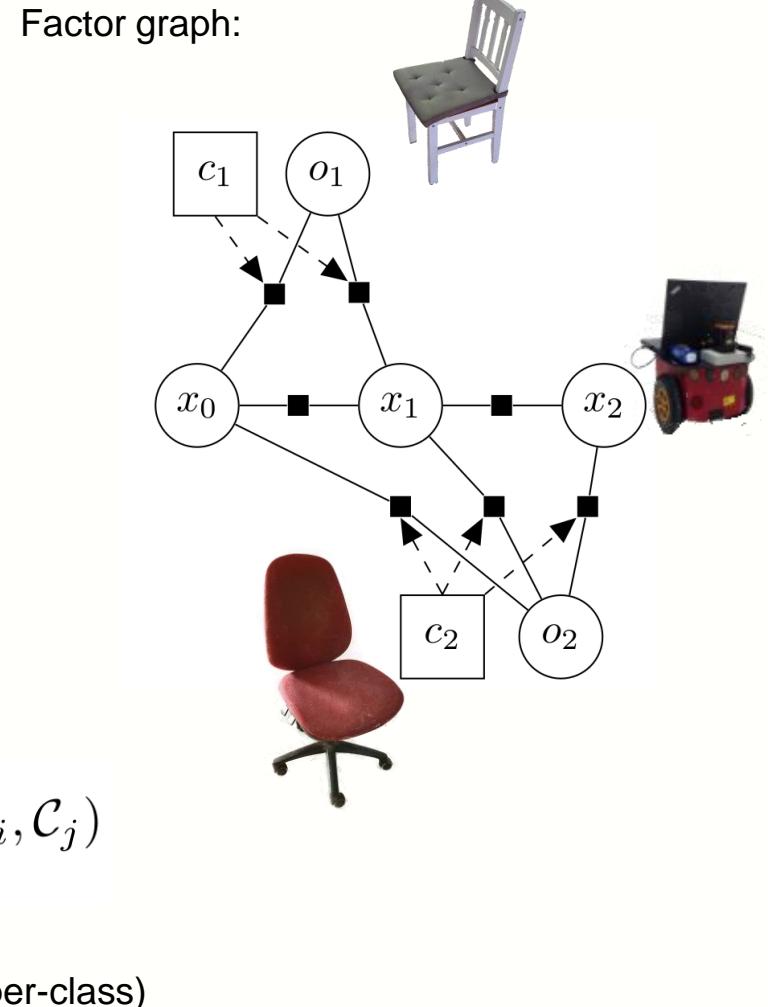
$$= \arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{O} \mid \mathcal{C}, \mathcal{H}_k) \cdot \mathbb{P}(\mathcal{C} \mid \mathcal{H}_k)$$

(continuous) hypothesis

hypothesis weight

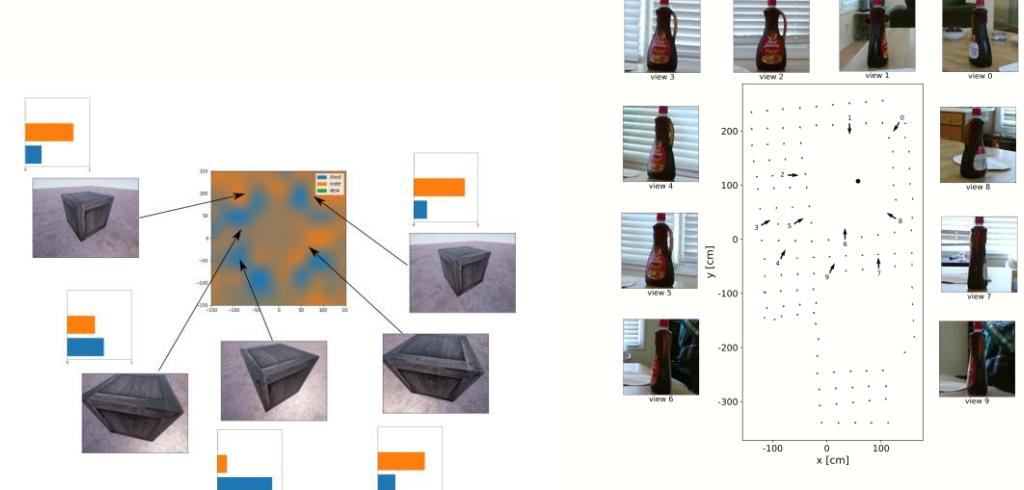
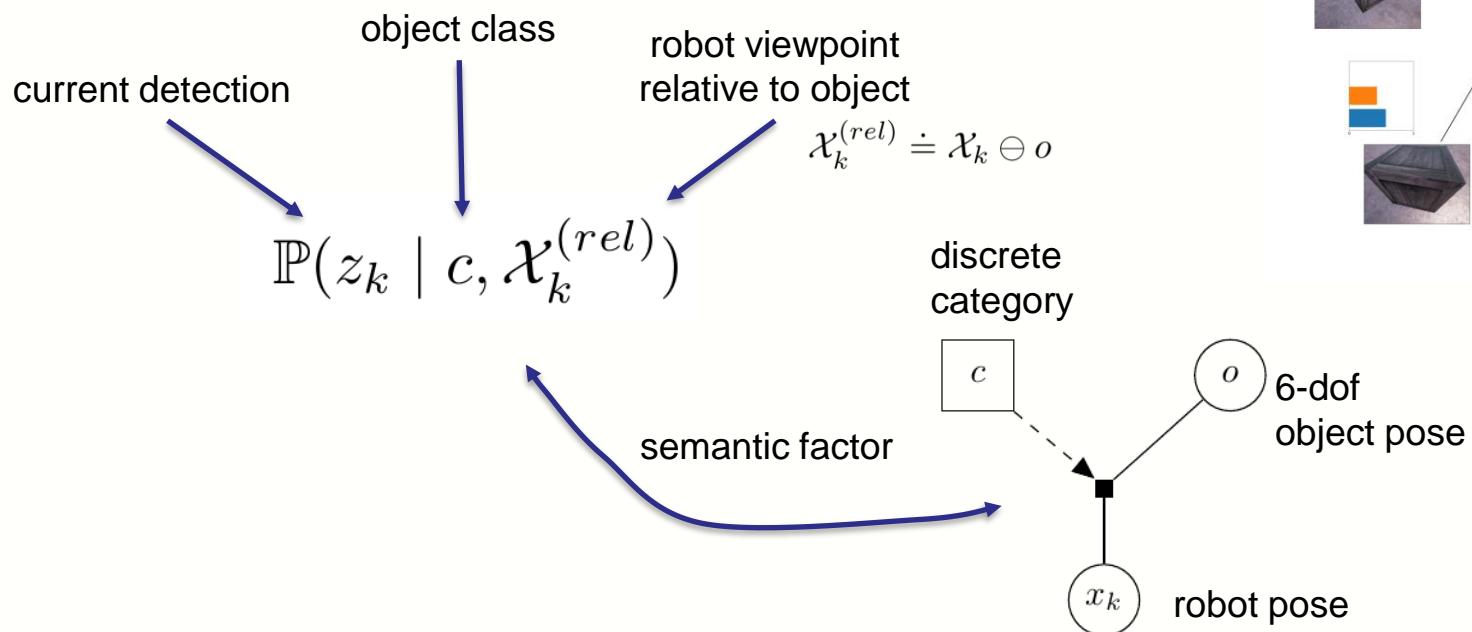
$$\mathbb{P}(\mathcal{X}_{0:k}, \mathcal{O} \mid \mathcal{C}, \mathcal{H}_k) = \eta \cdot \mathbb{P}(\mathcal{X}_0) \prod_i \mathbb{P}(\mathcal{X}_{i+1} \mid \mathcal{U}_i, \mathcal{X}_i) \prod_{j \in \mathcal{M}_i} \mathbb{P}(\mathcal{Z}_j \mid \mathcal{O}_j, \mathcal{X}_i, \mathcal{C}_j)$$

observation model (per-class)



Viewpoint-Dependent Models

- Viewpoint dependency \Rightarrow viewpoint-dependent models
 - Allow to couple semantics and geometry



Yuri Feldman and Vadim Indelman. "Spatially-dependent Bayesian semantic perception under model and localization uncertainty." Autonomous Robots (2020): 1-29.

Outline

Viewpoint-Dependent models for Semantic Perception under Uncertainty

- ✓ 1. The semantic perception problem (Object-Level SLAM)
- ✓ 2. Viewpoint-dependent semantic measurement models
- 3. Contributions:
 - I. Spatially-Dependent Classification under Model and Localization Uncertainty
 - II. Data Association-Aware Semantic Mapping and Localization
 - III. Semantic Perception with a Continuous Learned Representation

Spatially-Dependent Uncertainty-Aware Classification

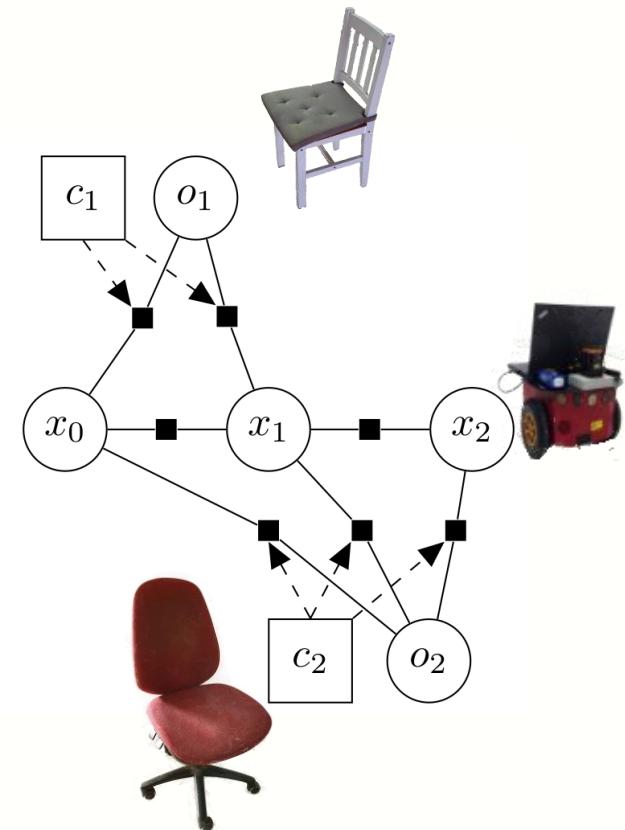
"Bayesian viewpoint-dependent robust classification under model and localization uncertainty", Feldman & Indelman 18' ICRA
"Spatially-dependent Bayesian semantic perception under model and localization uncertainty", Feldman & Indelman 20' ARJ

$$\arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{C}, \mathcal{O} | \mathcal{H}_k)$$

object categories (discrete!)
robot track
measurement and control history
object geometry

$$= \arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{O} | \mathcal{C}, \mathcal{H}_k) \cdot \mathbb{P}(\mathcal{C} | \mathcal{H}_k)$$

(continuous) hypothesis hypothesis weight



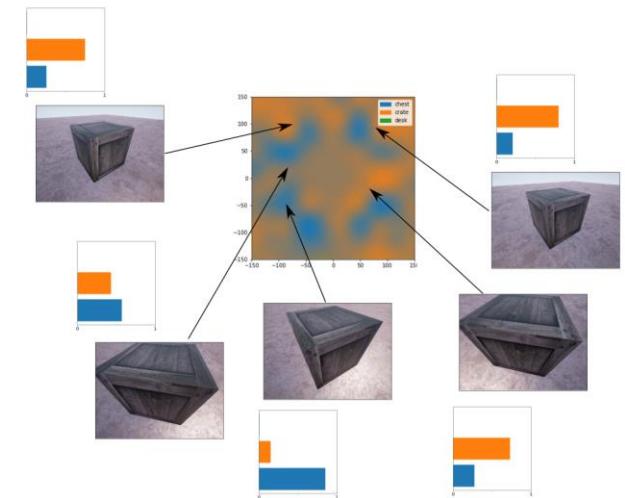
Semantic Measurements – Spatial Correlation

Semantic measurements spatially correlated \Rightarrow not i.i.d.

Would like to model the joint likelihood

$$\mathbb{P}(\mathcal{S}_{0:k} \mid c, X_{0:k}^{(rel)})$$

object class
classifier responses track of relative poses



Can be done by fitting a Gaussian Process to classifier responses (offline training step)

$$s = f_c(x^{(rel)}) + \epsilon$$

$$f_c(x^{(rel)}) \sim \mathcal{GP} \left(\mu_c(x^{(rel)}), k_c(\cdot, \cdot) \right)$$

Semantic Measurements – Model Uncertainty

Problem: DNN output away from training set unstable

Solution (Gal & Ghahramani, 16' and others):

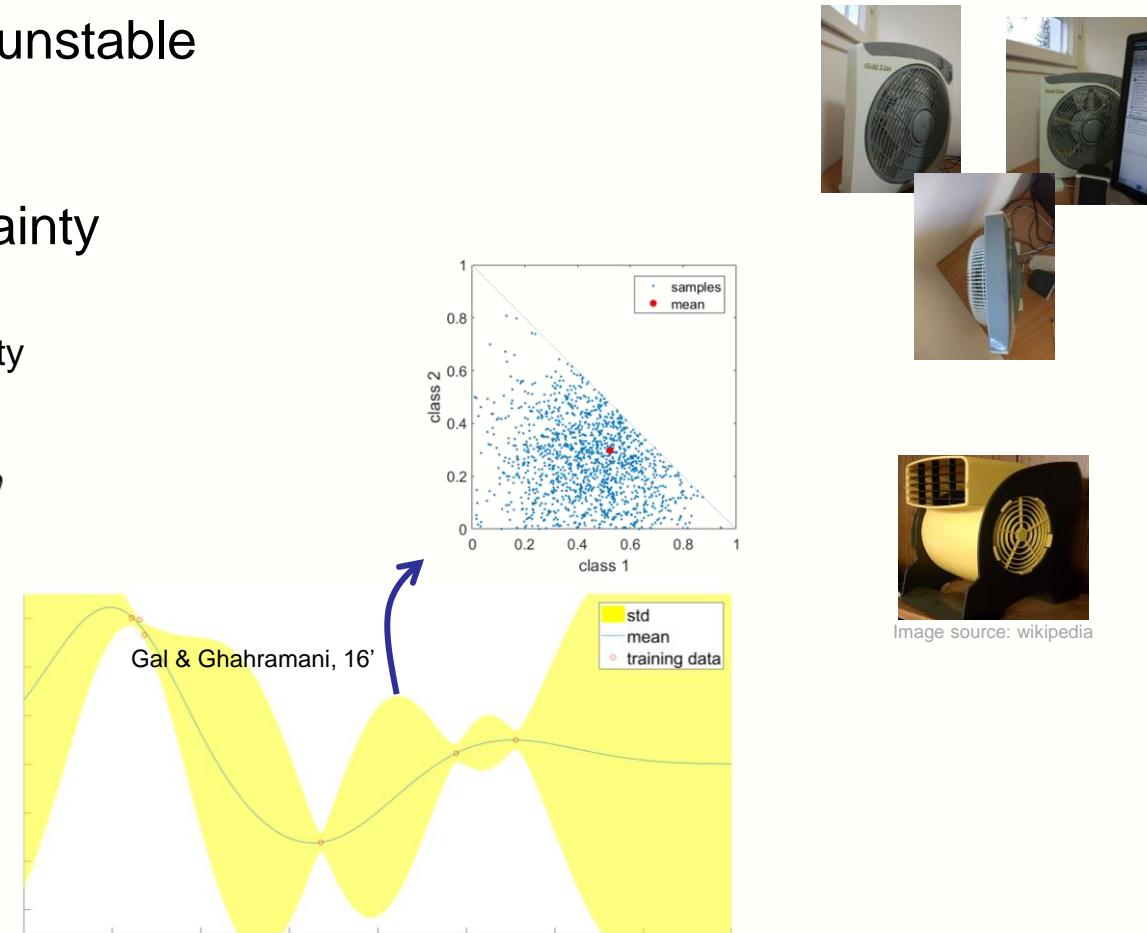
- ❖ marginalize over weight (model) uncertainty

Formally network output sample model uncertainty

$$\mathbb{P}(c | \mathcal{D}) = \int \mathbb{P}(c | \mathcal{W}) \cdot \mathbb{P}(\mathcal{W} | \mathcal{D}) d\mathcal{W}$$

G&G 16': $\mathbb{P}(\mathcal{W} | \mathcal{D}) \approx \text{Bernoulli}\left(\frac{1}{2}\right) \cdot \hat{\mathcal{W}}$

Marginalization approximated via importance sampling.



Approach – Uncertainty-Aware Classification

$$\mathbb{P}(c \mid \mathcal{H}_k) = \int_{\mathcal{X}_{0:k}, o} \underbrace{\mathbb{P}(c \mid \mathcal{X}_{0:k}, o, \mathcal{H}_k)}_{(a)} \underbrace{\mathbb{P}(\mathcal{X}_{0:k}, o \mid \mathcal{H}_k)}_{(b)} d\mathcal{X}_{0:k} do$$

Marginalize over last classification Marginalize over Landmarks

$$\frac{1}{n_k} \sum_{s_k \in \mathcal{S}_k} \mathbb{P}(c \mid s_k, H_k \setminus \{z_k\})$$

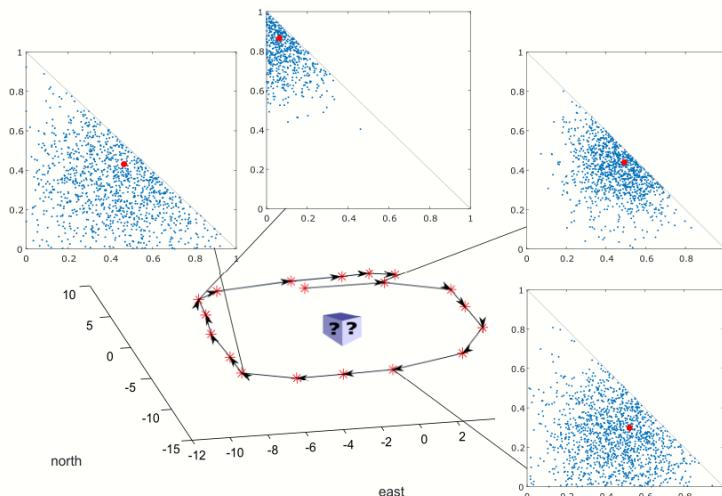
Bayes

$$\frac{\mathbb{P}(s_k \mid c, H_k \setminus \{z_k\}) \cdot \mathbb{P}(c \mid H_{k-1})}{\sum_{c \in \mathcal{C}} \mathbb{P}(s_k \mid c, H_k \setminus \{z_k\}) \cdot \mathbb{P}(c \mid H_{k-1})}$$

Marginalize over past classifications

$$\int_{\mathcal{S}_{0:k-1}} \mathbb{P}(s_k \mid c, \mathcal{S}_{0:k-1}, \mathcal{X}_{0:k}^{(rel)}) \cdot \prod_{i=0}^{k-1} \mathbb{P}(s_i \mid z_i) d\mathcal{S}_{0:k-1}$$

Class model Model uncertainty



Results

Evaluated in synthetic simulation, 3D simulation, real-world data (BigBIRD, AVD).

Single object classification from measurements over a track. Localization uncertain but estimate available.

Evaluation criteria:

1. Probability of ground-truth class (higher is better) $\mathbb{P}(c^{GT} \mid \mathcal{H})$
 2. Most-likely-to-ground-truth ratio (lower is better)
 \Rightarrow sensitive to confident misclassifications
- $$MGR \doteq \frac{\arg \max_c \mathbb{P}(c \mid \mathcal{H})}{\mathbb{P}(c^{GT} \mid \mathcal{H})}$$

Baselines:

- ❖ “Model Based” Teacy et al. 15’ AAMAS – GP class model, assumes known localization
- ❖ Naïve Bayes (synthetic simulation) – directly fuses class predictions, no class model

Results – Synthetic Simulation

Statistics for several hand-specified scenarios over realizations of simulated classification.
3 candidate classes with hand-specified GP models, measurements from ground truth GP.

Model uncertainty benchmark

- Simulated classification is randomly offset at each step.
Our method is input with uncertainty.

Localization uncertainty benchmark

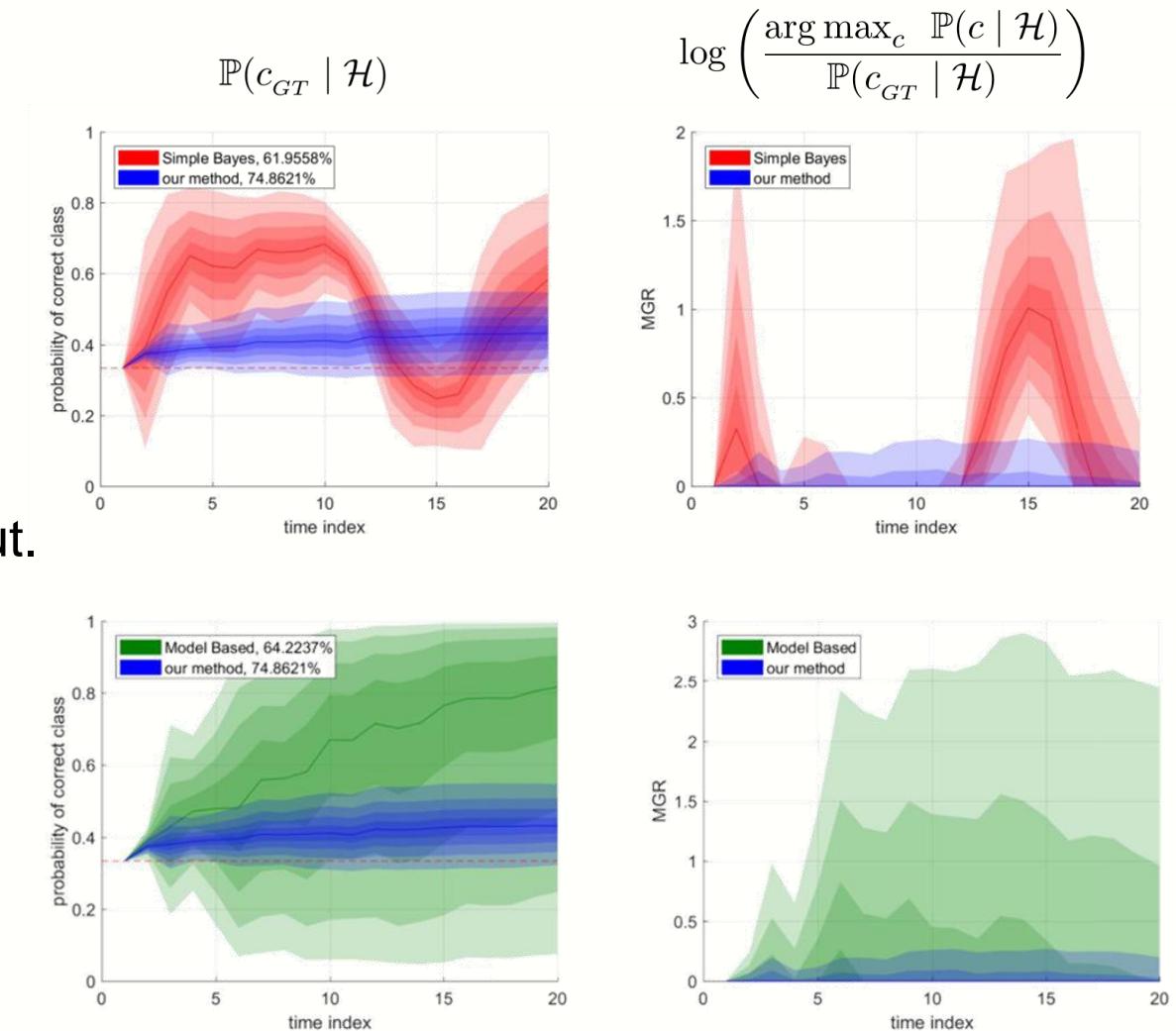
- Localization biased in a way that creates aliasing.
Our method is input with uncertainty.

Results – Synthetic Simulation – Model Uncertainty

Color patches are equal-step (10%) percentiles.
Saturated line is median.
Legend lists % of steps with GT class most likely.

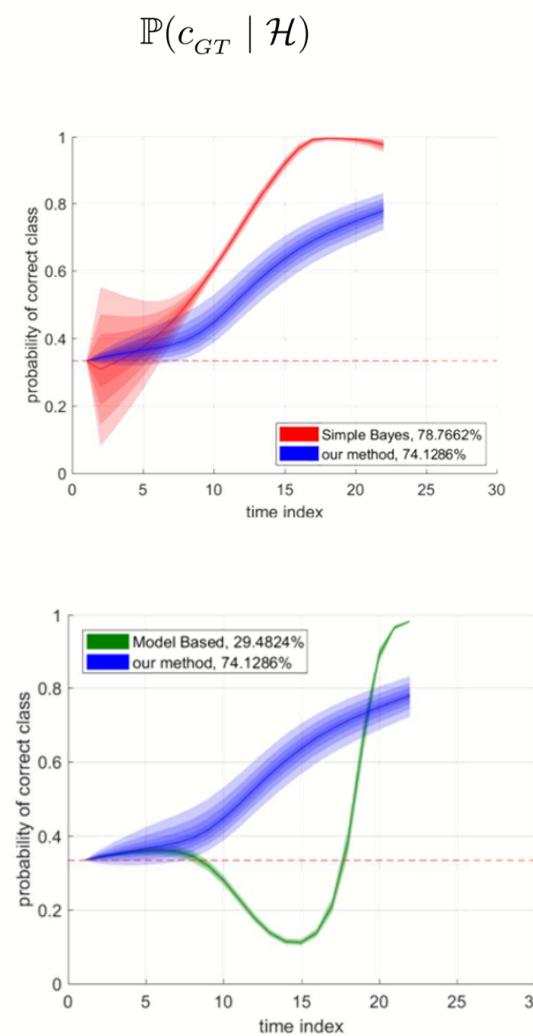
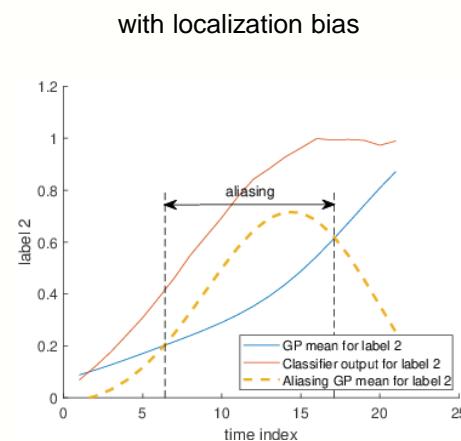
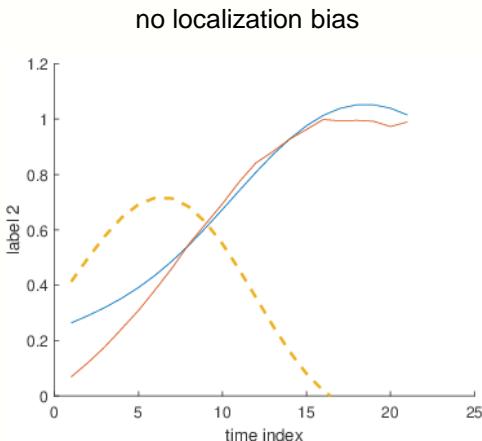
Results:

- ❖ Naïve Bayes arbitrarily off due to erroneous input.
- ❖ Model Based gives high scores to GT, but misclassifies in nearly 30% of the cases (MGR).
- ❖ Our method gracefully accumulates information, dense performance percentiles.

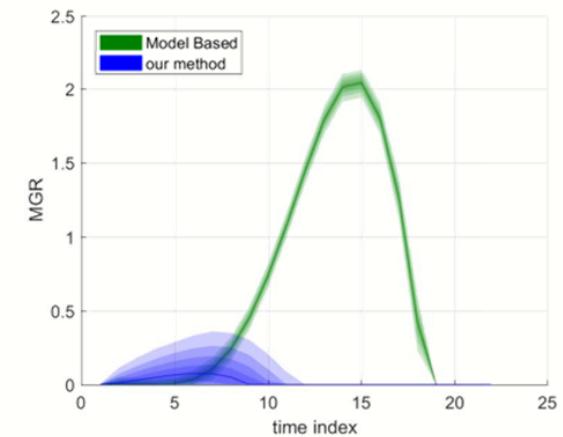
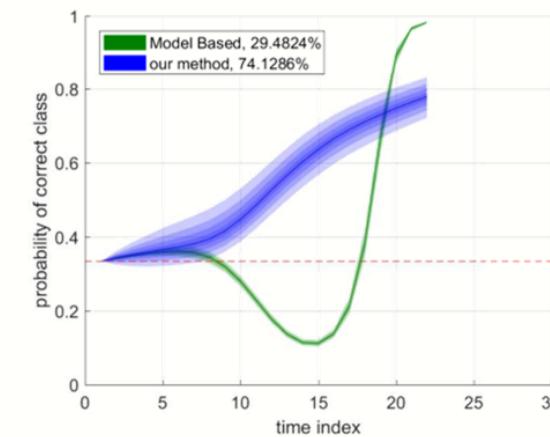
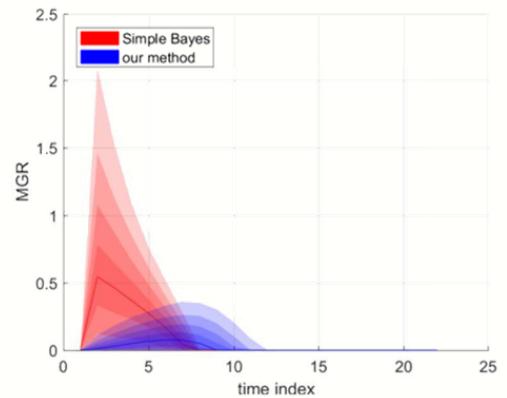


Results – Synthetic Simulation – Localization Uncertainty

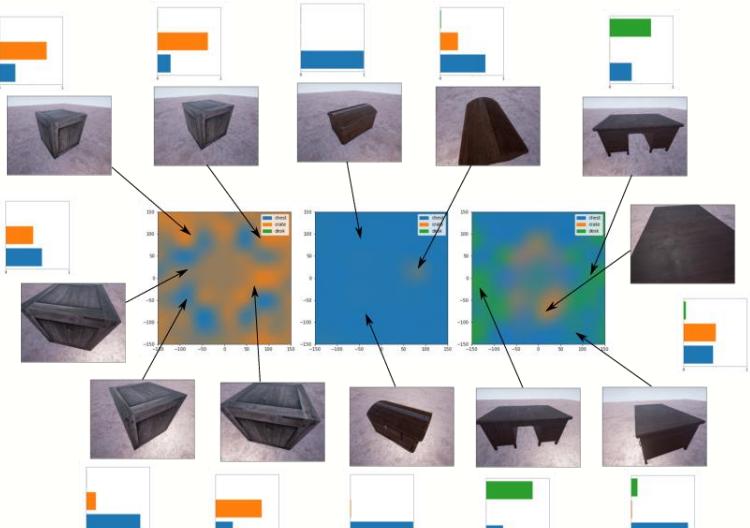
- ❖ Model Based arbitrarily off due to aliasing.
- ❖ Naïve Bayes unaffected (classification measurements are correct).
- ❖ Our method gracefully accumulates information.



$$\log \left(\frac{\arg \max_c \mathbb{P}(c | \mathcal{H})}{\mathbb{P}(c_{GT} | \mathcal{H})} \right)$$

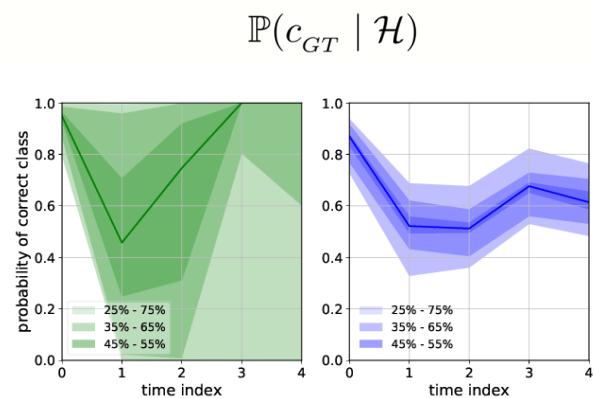


Learned GP models

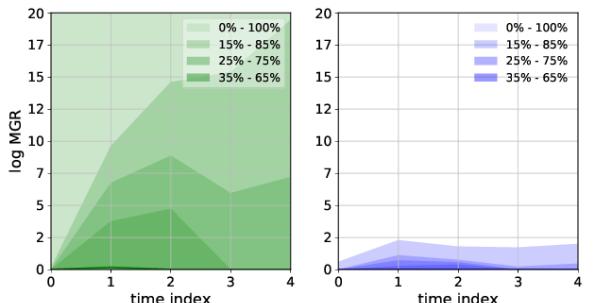


Results – 3D Simulation

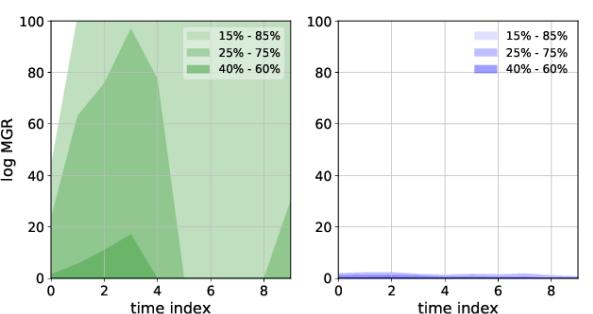
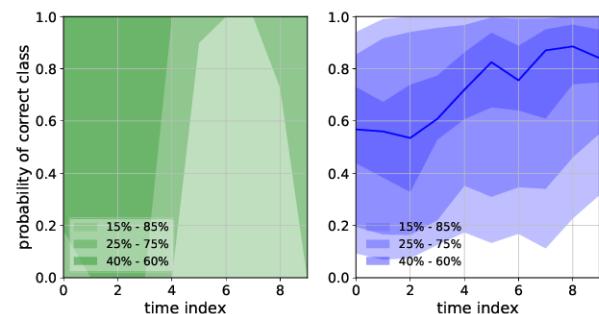
Localization uncertainty:



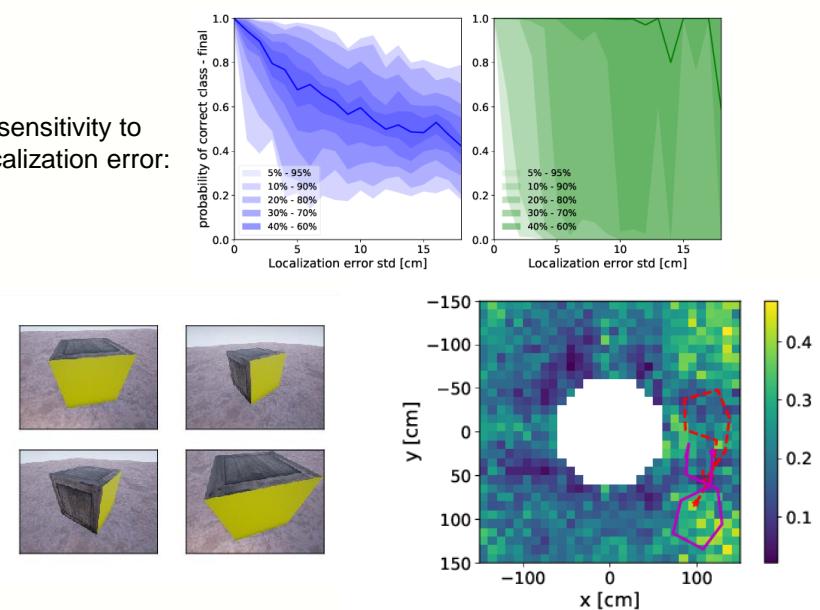
$$\log \left(\frac{\arg \max_c \mathbb{P}(c | \mathcal{H})}{\mathbb{P}(c_{GT} | \mathcal{H})} \right)$$



Model uncertainty:

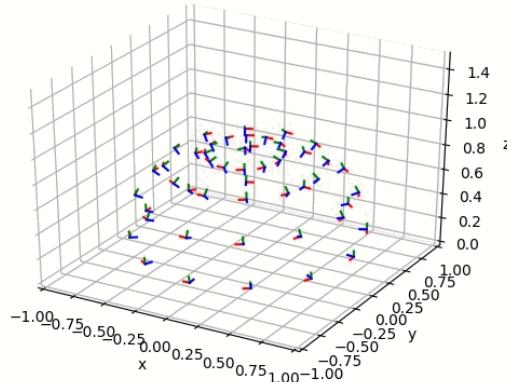
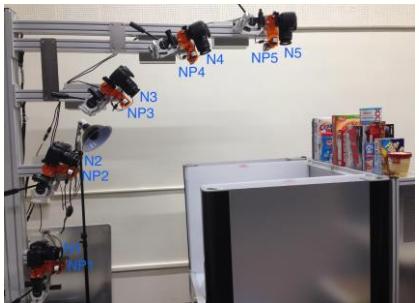


sensitivity to
localization error:

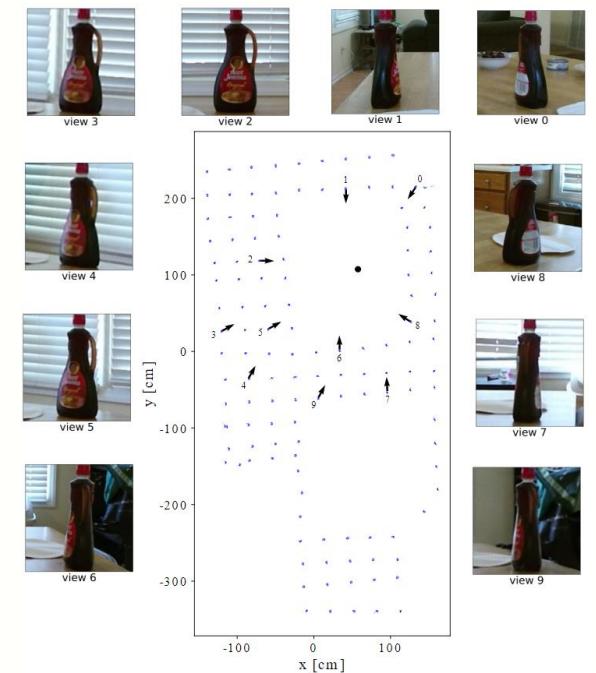
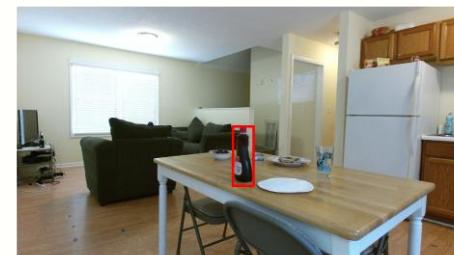


Results – Real Data

BigBIRD
(Big Berkeley Instance Recognition)
⇒ 125 objects
120 views / object

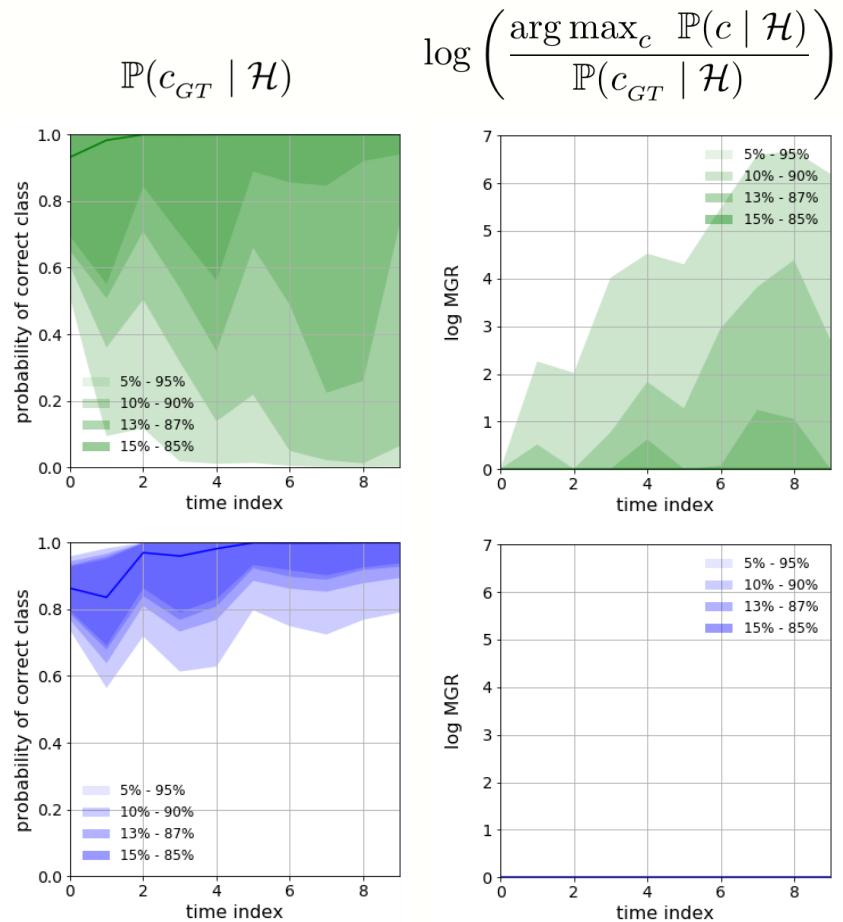


AVD (Active Vision Dataset)

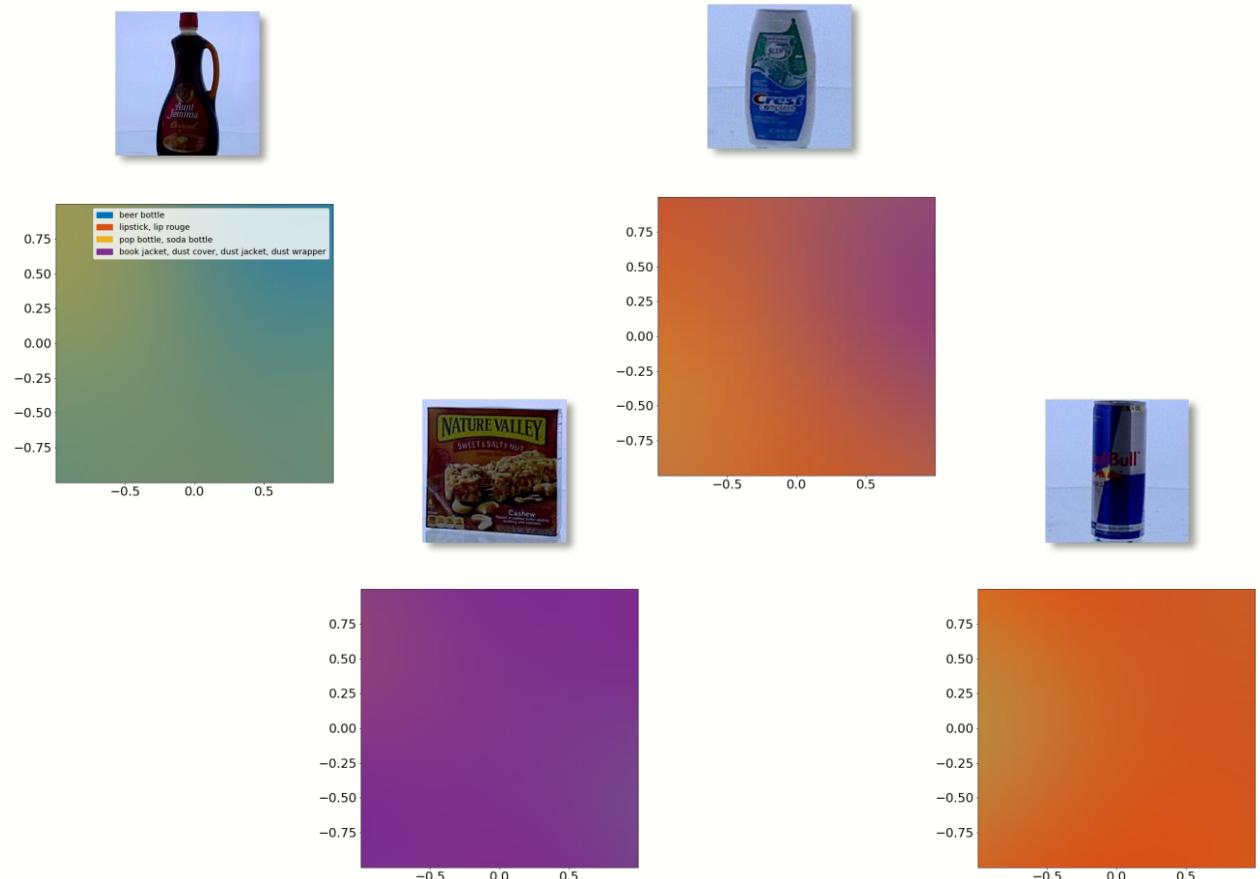


Results – Real Data

Localization uncertainty



Example objects and learned GPs



Outline

Viewpoint-Dependent models for Semantic Perception under Uncertainty

- ✓ 1. The semantic perception problem (Object-Level SLAM)
- ✓ 2. Viewpoint-dependent semantic measurement models
- 3. Contributions:
 - I. Classification under Model and Localization Uncertainty
 - II. Data Association-Aware Semantic Mapping and Localization
 - III. Semantic Perception with a Continuous Learned Representation

DA - Aware Semantic Mapping and Localization

“Data association aware semantic mapping and localization via a viewpoint-dependent classifier model”, Tchuiev, Feldman and Indelman 19’ IROS

Up until now: assumed data association solved.

Consider data association:

Consider data association:

$$\arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}, \beta} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{C}, \mathcal{O}, \beta_{0:k} | \mathcal{H}_k)$$

data association per measurement (discrete)

measurement and control history

robot track (continuous)

object categories (discrete)

object geometry (continuous)

$$= \arg \max_{\mathcal{X}, \mathcal{C}, \mathcal{O}, \beta} \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{O} | \mathcal{C}, \beta_{0:k}) \cdot \mathbb{P}(\mathcal{C}, \beta_{0:k} | \mathcal{H}_k)$$

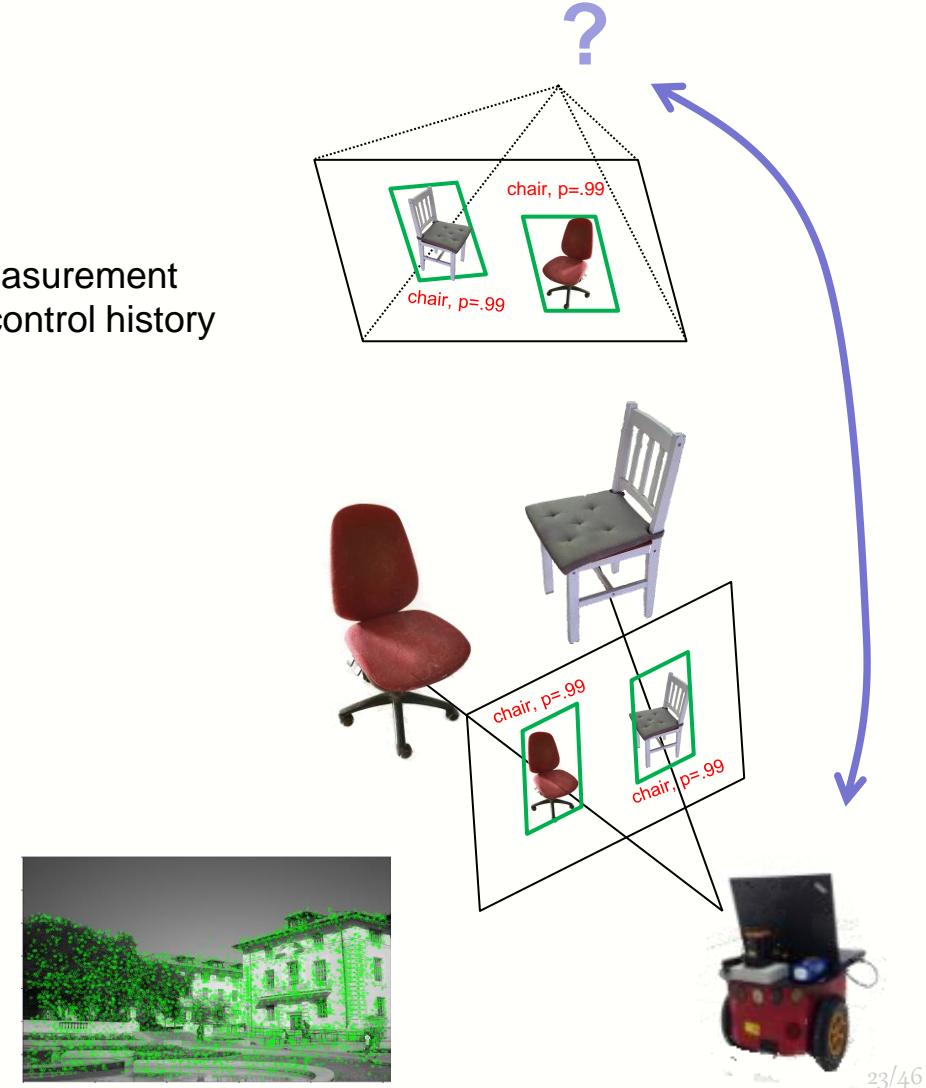
(continuous) hypothesis hypothesis weight

$c \in \{1, \dots, N\}$

...which is still better than ...



$$|\{(c_1, \dots, c_m, \beta_1, \dots, \beta_k)\}| = N^M \cdot m^k$$



DA - Aware Semantic Mapping and Localization

Contributions:

- ❖ DA-Aware Semantic Mapping through maintaining a full hybrid belief.
 - ⇒ Maintain all plausible hypotheses until disambiguation is possible.
Tractable in practice (subject to pruning).
- ❖ Viewpoint-Dependent model aids DA disambiguation by coupling between geometry and semantics.
- ❖ Approach operates with rich semantic feature vectors, not limited to most-likely-class measurements.

Approach - DA - Aware Semantic SLAM

Denote $b[\mathcal{X}_{0:k}, \mathcal{O}]_{\beta_{1:k}}^{\mathcal{C}} \doteq \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{O} \mid \mathcal{C}, \beta_{1:k})$ (continuous) hypothesis $w_{\beta_{1:k}}^{\mathcal{C}} \doteq \mathbb{P}(\mathcal{C}, \beta_{1:k} \mid \mathcal{H}_k)$ hypothesis weight

hypothesis propagation $b[\mathcal{X}_{0:k}, \mathcal{O}]_{\beta_{1:k}}^{\mathcal{C}} \propto b[\mathcal{X}_{0:k-1}, \mathcal{O}]_{\beta_{1:k-1}}^{\mathcal{C}} \cdot \mathbb{P}(\mathcal{X}_k \mid \mathcal{X}_{k-1}, \mathcal{A}_{k-1}) \cdot \mathbb{P}(\mathcal{Z}_k \mid \mathcal{X}_k, \mathcal{O}_{\beta_k}, \mathcal{C})$ motion model viewpoint-dependent model

weight update $w_{\beta_{1:k}}^{\mathcal{C}} \propto w_{\beta_{1:k-1}}^{\mathcal{C}} \int_{\mathcal{X}, \mathcal{O}} \mathbb{P}(\beta_k \mid \mathcal{X}_k, \mathcal{O}_{\beta_k}) \cdot b[\mathcal{X}_{0:k}, \mathcal{O}]_{\beta_{1:k}}^{\mathcal{C}} d\mathcal{X} d\mathcal{O}$ association probability

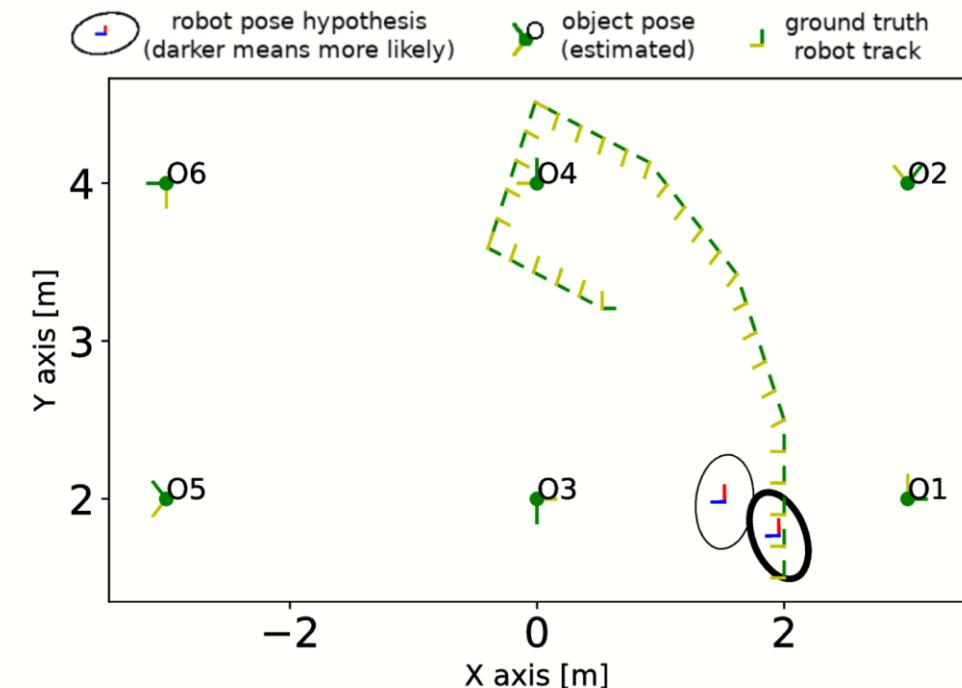
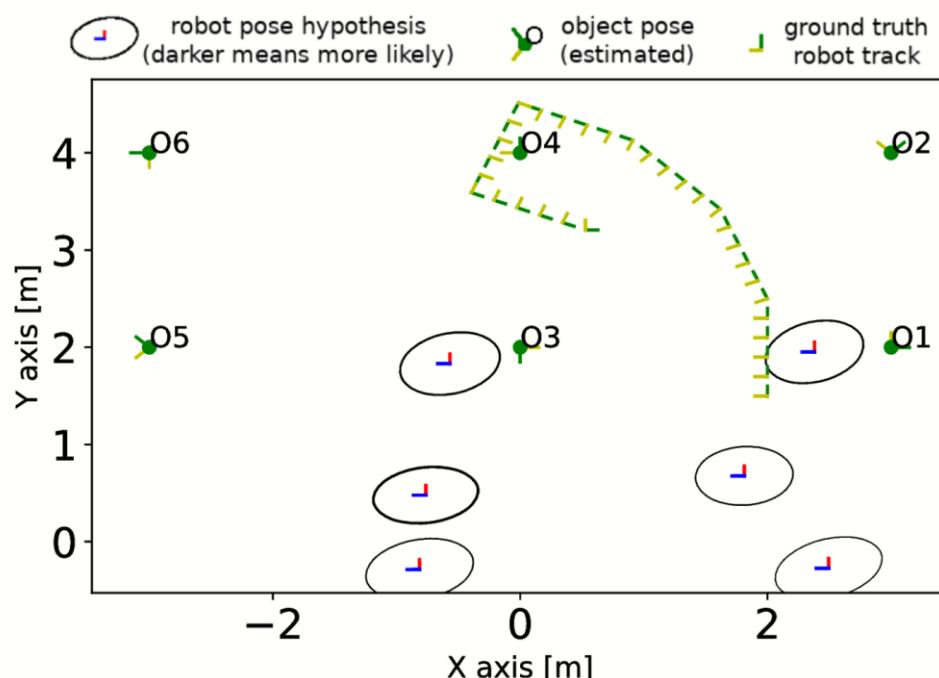
Weights that fall below a threshold are pruned

Results

Simulated environment:

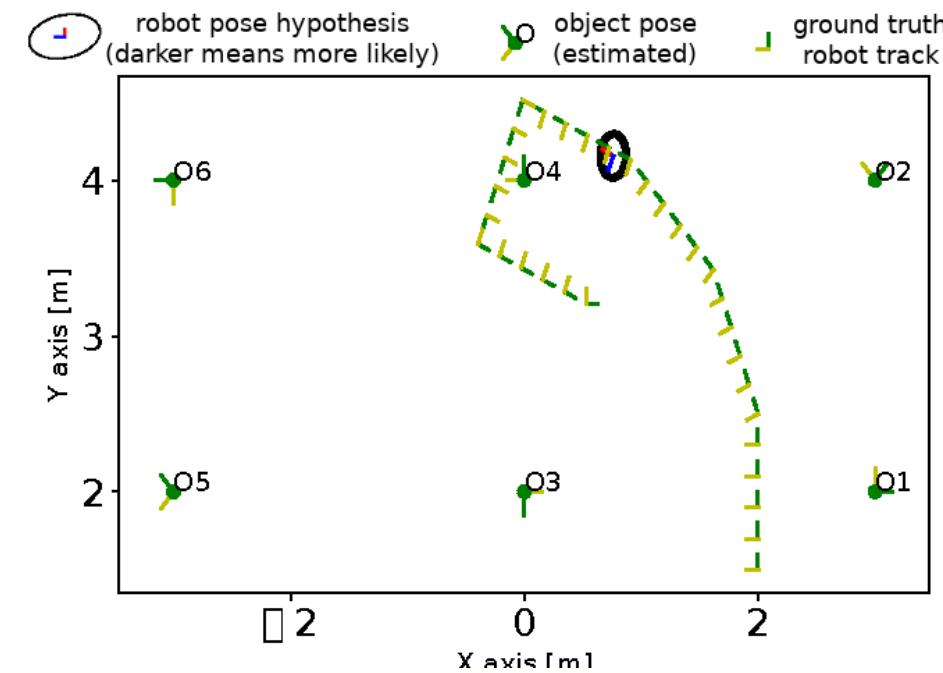
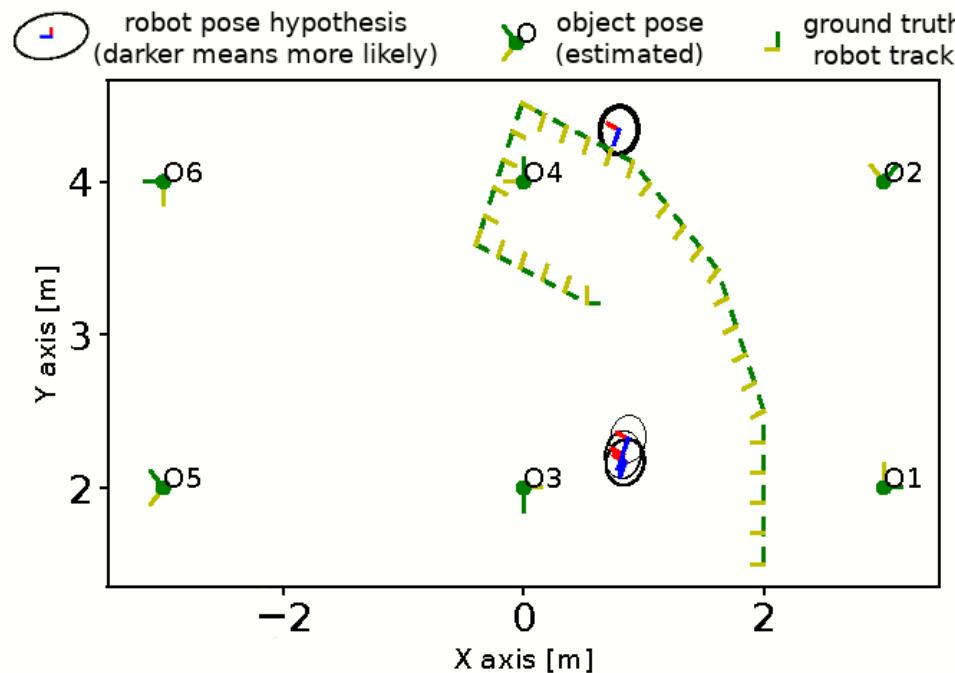
- ❖ 6 identical objects
- ❖ uninformative robot pose prior
- ❖ 2 candidate classes with synthetic measurement models

Time $k = 1$, without (left) and with (right) classifier model.



Results

Time $k = 15$ without (left) and with (right) classifier model

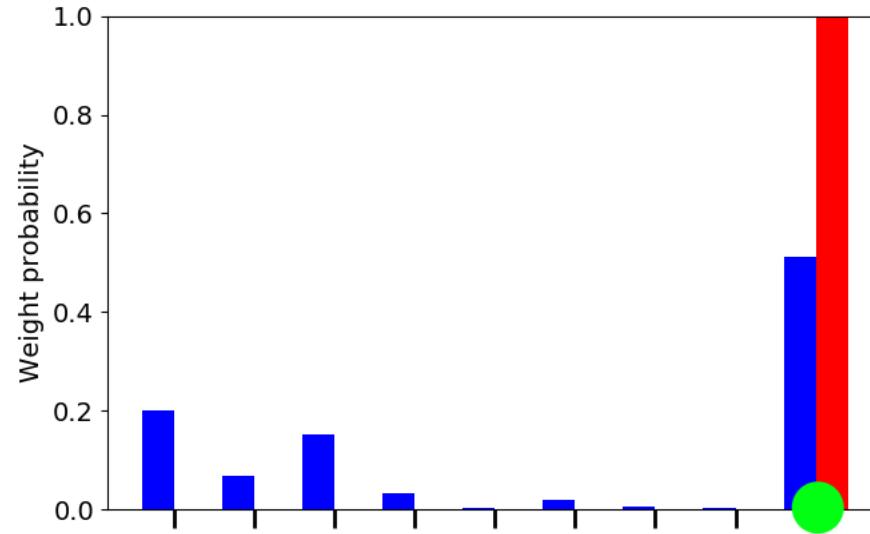
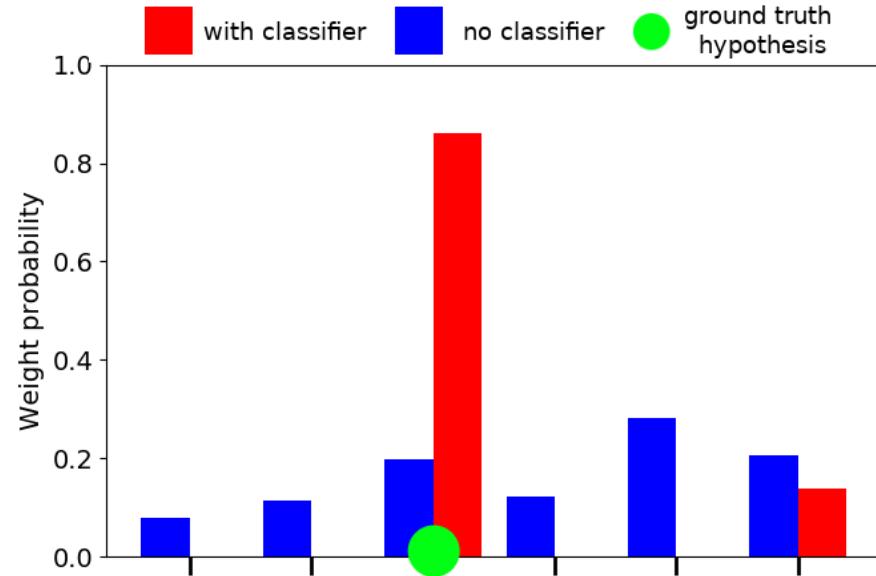


With classifier:

- ❖ Fewer hypotheses
- ❖ More accurate localization

Results

Hypothesis weight comparison, times $k = 1$ (left) and $k = 15$ (right)



With classifier:

- ❖ Fewer hypotheses
- ❖ Stronger disambiguation

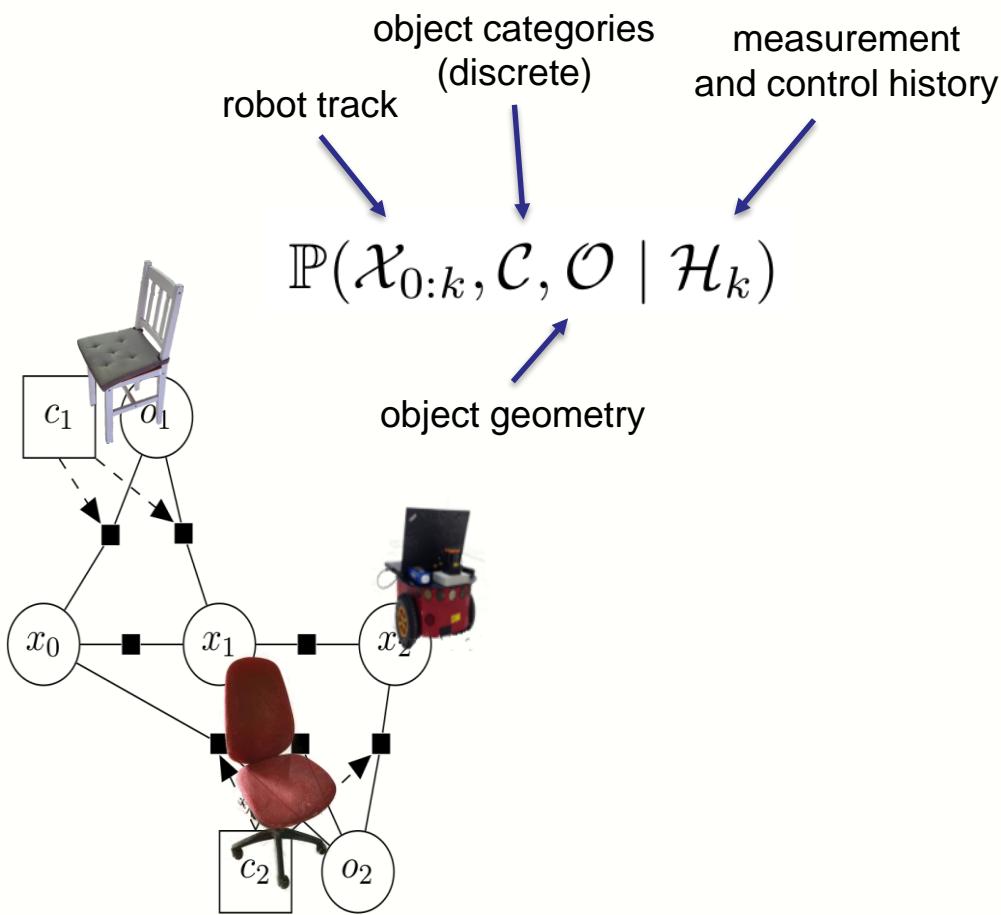
Outline

Viewpoint-Dependent models for Semantic Perception under Uncertainty

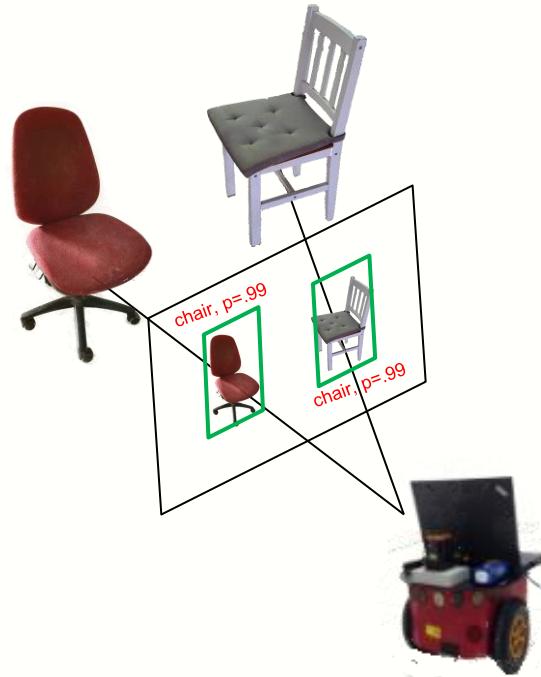
- ✓ 1. The semantic perception problem (Object-Level SLAM)
- ✓ 2. Viewpoint-dependent semantic measurement models
- 3. Contributions:
 - I. Classification under Model and Localization Uncertainty
 - II. Data Association-Aware Semantic Mapping and Localization
 - III. Semantic Perception with a Continuous Learned Representation

Semantic Perception with a Continuous Learned Representation

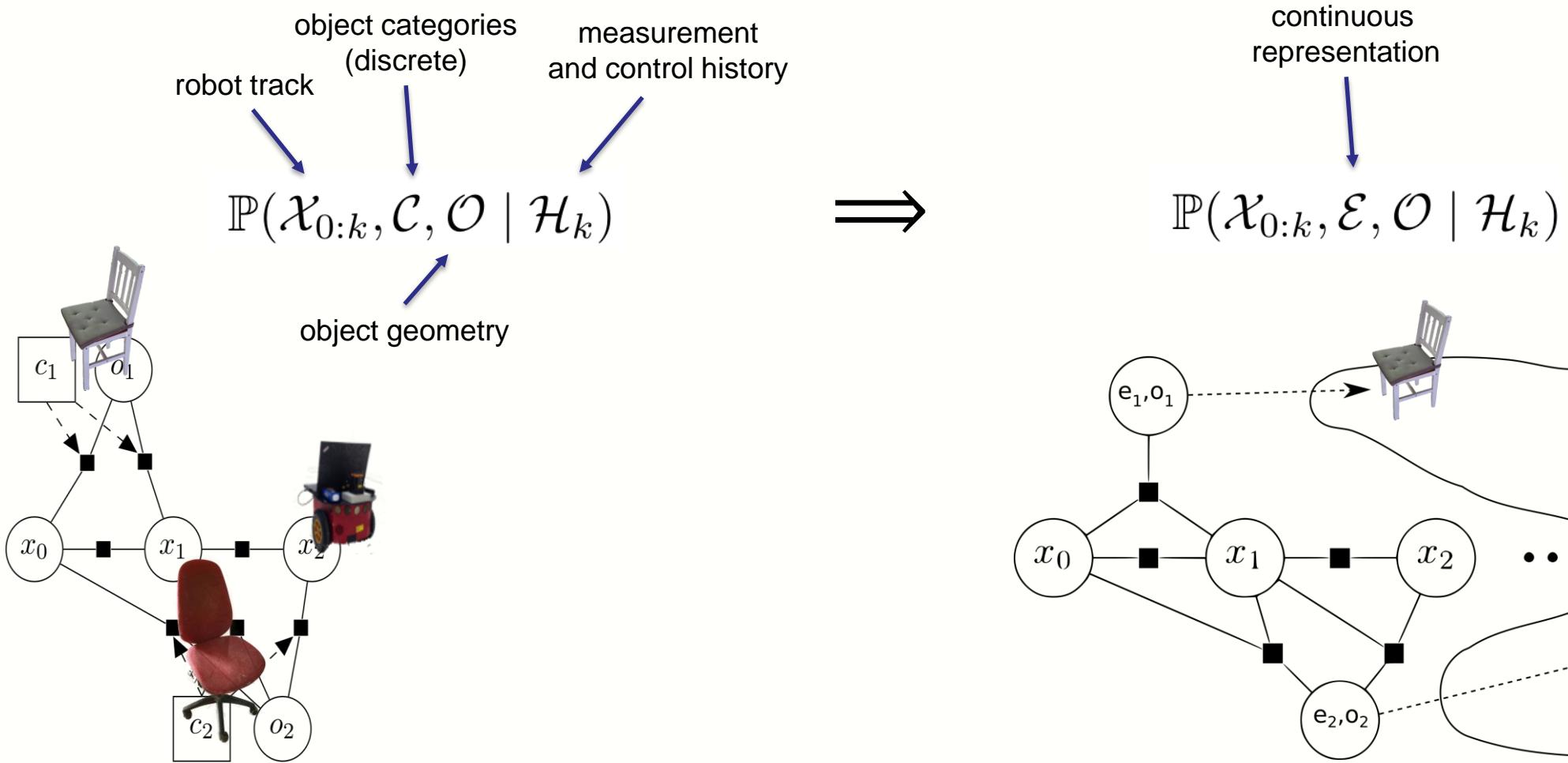
An initial version presented as “Towards Self-Supervised Semantic Representation with a Viewpoint-Dependent Observation Model” in proceedings of Workshop on Self-Supervised Robot Learning, in conjunction with RSS, July 2020



- ⇒ Requires maintaining hypotheses (inefficient)
- Requires per-class models
- Limited granularity of semantic representation

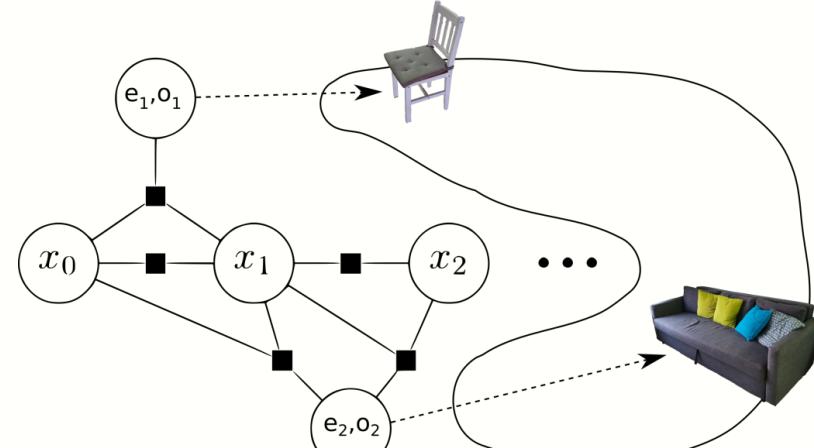


Semantic Perception with a Continuous Learned Representation



Continuous Learned Representation – Inference

$$\begin{aligned}
 & \text{maintain joint belief: } b_k \doteq \mathbb{P}(\mathcal{X}_{0:k}, \mathcal{E}, \mathcal{O} \mid \mathcal{H}_k) \\
 & \quad \text{robot track} \quad \text{continuous representation} \quad \text{measurement and control history} \\
 & \quad \text{object geometry} \\
 & = \eta \cdot \underbrace{\mathbb{P}(z_k \mid \mathcal{X}_{0:k}, \mathcal{E}, \mathcal{O}, \mathcal{H}_k \setminus \{z_k\})}_{\text{semantic observation model}} \\
 & \quad \cdot \underbrace{\mathbb{P}(\mathcal{X}_k \mid \mathcal{X}_{0:k-1}, \mathcal{A}_{k-1}) \cdot b_{k-1}}_{\text{motion model}}
 \end{aligned}$$



- ⇒ Need a single semantic observation model (as opposed to per-class previously)
- Continuous inference
- No discretization on semantic representation

Fitting the Viewpoint-Dependent Model – Take 1

Assume a Gaussian viewpoint-dependent model conditioned on continuous representation:

$$\mathbb{P}(\mathcal{Z}_k \mid \mathcal{X}_{0:k}, \mathcal{E}, \mathcal{O}, \mathcal{H}_k \setminus \{z_k\}) \doteq \mathbb{P}(\mathcal{Z}_k \mid \mathcal{X}_k^{(rel)}, \mathcal{E}) \\ \doteq \mathcal{N}(\mathcal{Z}_k; \mu_\theta(\mathcal{X}_k^{(rel)}, \mathcal{E}), \Sigma)$$

frame
relative pose to object semantic description

Use Maximum-a-Posteriori to fit \mathcal{E}, θ

$$\arg \max_{\theta, \mathcal{E}_{1:n}} \mathbb{P}(\mathcal{Z}_{0:k}, \mathcal{E}_{1:n} \mid \mathcal{X}_{0:k}^{(rel)}, \beta_{0:k})$$

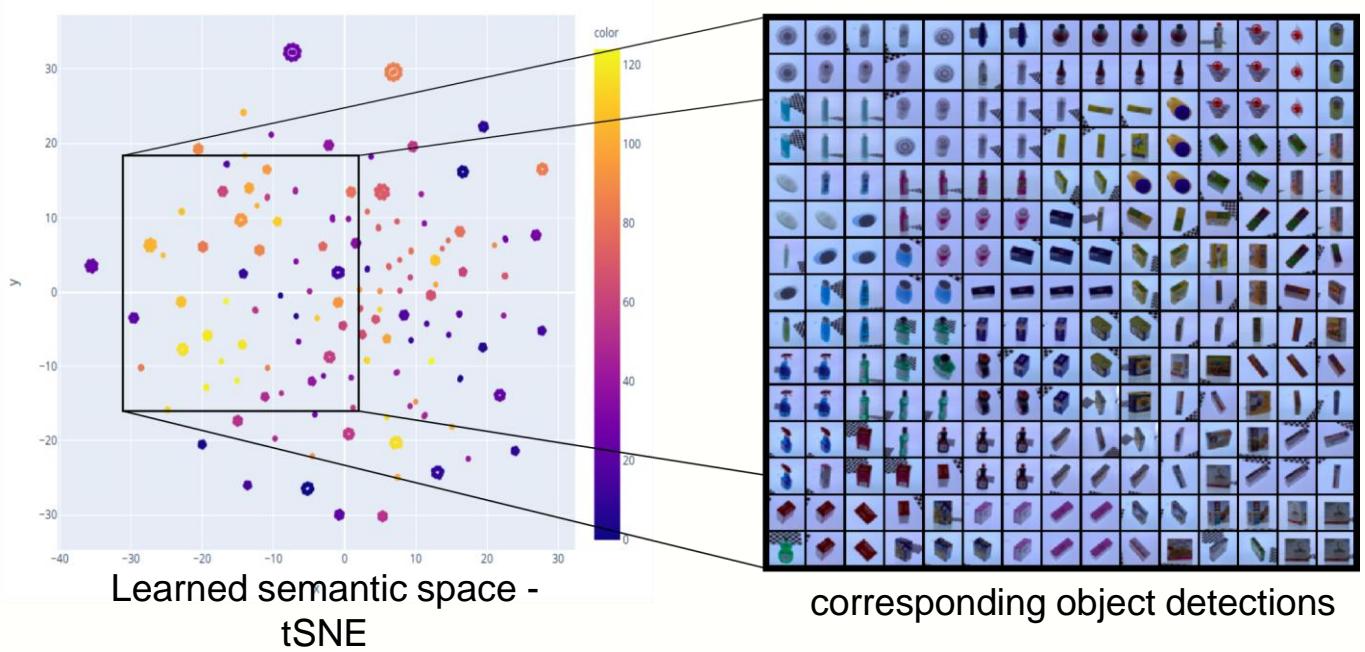
observations representation variables for n objects corresponding relative poses data association

$$= \arg \max_{\theta, \mathcal{E}_{1:n}} \sum \log \underbrace{\mathbb{P}_\theta(\mathcal{Z}_i \mid \mathcal{E}_{\beta_i}, \mathcal{X}_i^{(rel)})}_{\text{viewpoint - dependent observation model}} + \underbrace{\log \mathbb{P}(\mathcal{E}_{1:n})}_{\text{prior on representation vectors}}$$

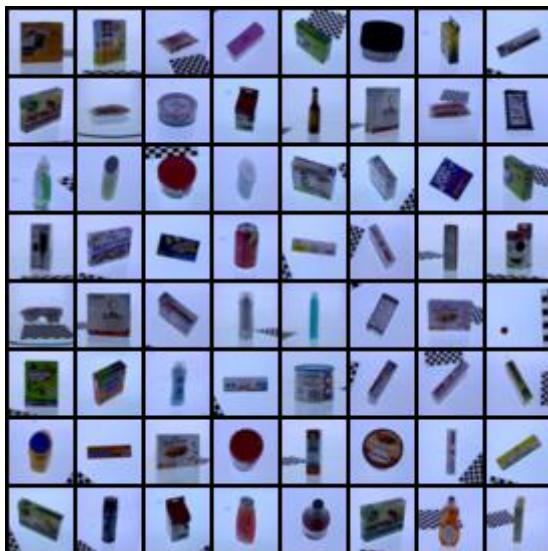
Take 1 - Results for Fitting $\mathbb{P}_\theta \left(\mathcal{Z} \mid \mathcal{E}, \mathcal{X}^{(rel)} \right)$



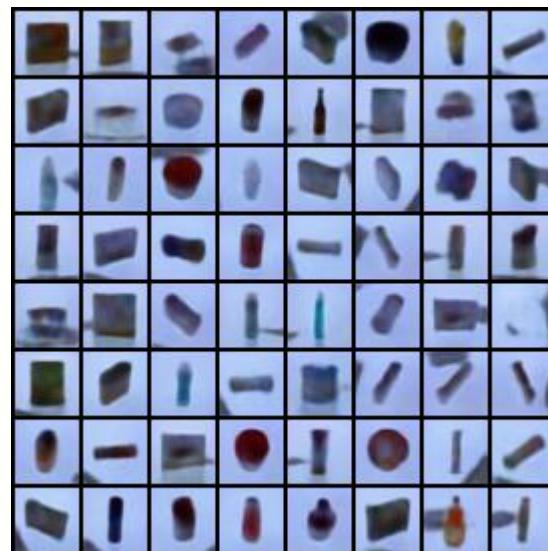
example images



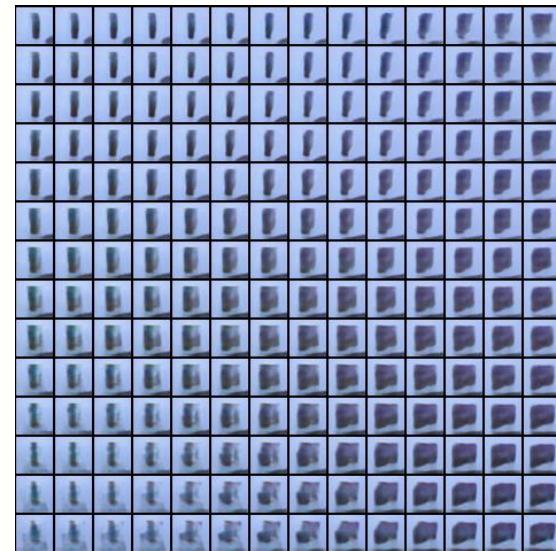
Take 1 - Results for Fitting $\mathbb{P}_\theta \left(\mathcal{Z} \mid \mathcal{E}, \mathcal{X}^{(rel)} \right)$



example images



mean predictions
(at ground truth)

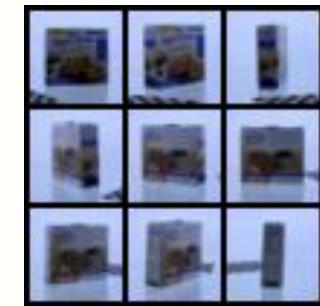
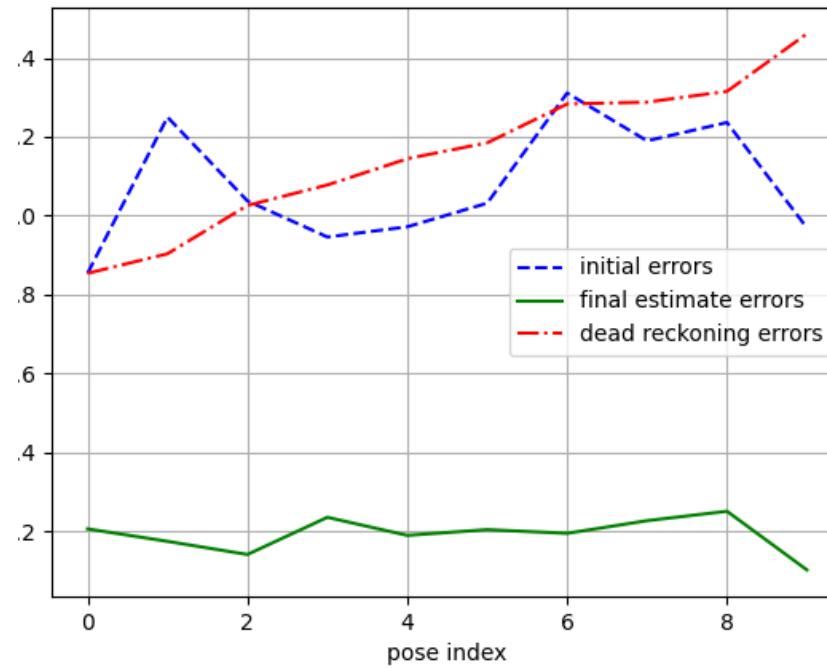
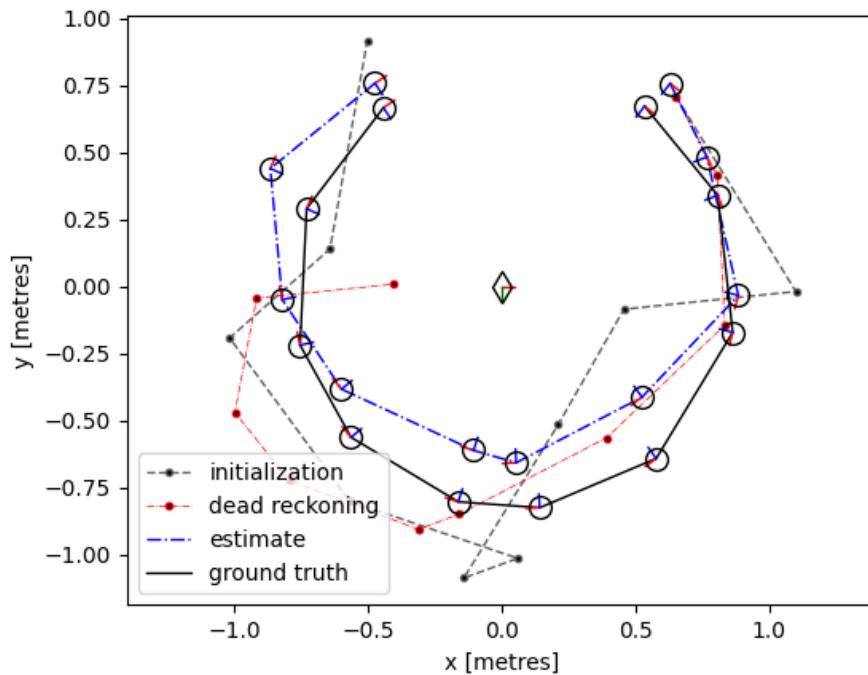


predictions for varying X
(around ground truth)

Take1 - Inference Using the Model

Simulation: use frames from viewpoints along a simulated track

- ❖ Also using odometry



example track frames

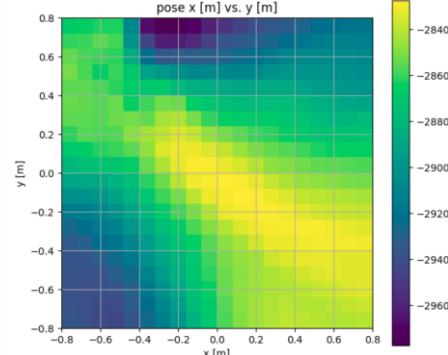
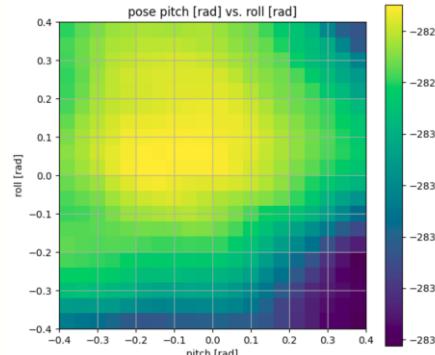
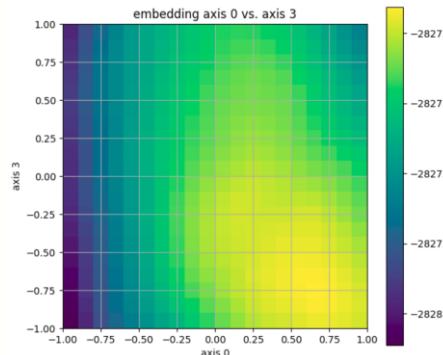
Take 1 - Limitations

- ❖ Factor is huge (32×32 frame \Rightarrow 1024 Jacobian rows / keyframe!)
- ❖ Model expressiveness?
- ❖ Maximum likelihood does not provide sufficient gradients for optimization

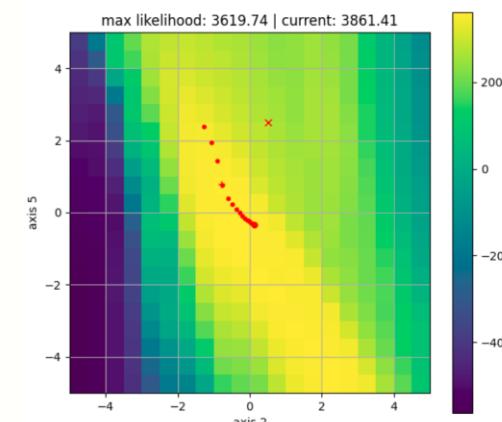
$$\arg \max_{\theta, \mathcal{E}_{1:n}} \sum \log \mathbb{P}_\theta \left(\mathcal{Z}_i \mid \mathcal{E}_{\beta_i}, \mathcal{X}_i^{(rel)} \right) + \log \mathbb{P}(\mathcal{E}_{1:n})$$

only constraints at training examples!

Values of likelihood $\mathbb{P}_\theta \left(\mathcal{Z} \mid \mathcal{E}, \mathcal{X}^{(rel)} \right)$ around ground truth point:



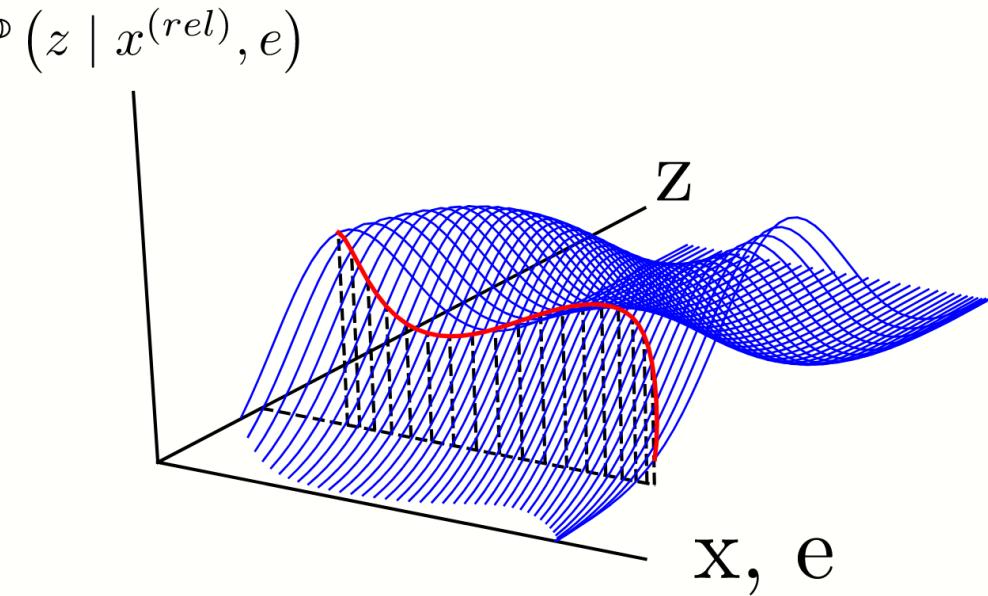
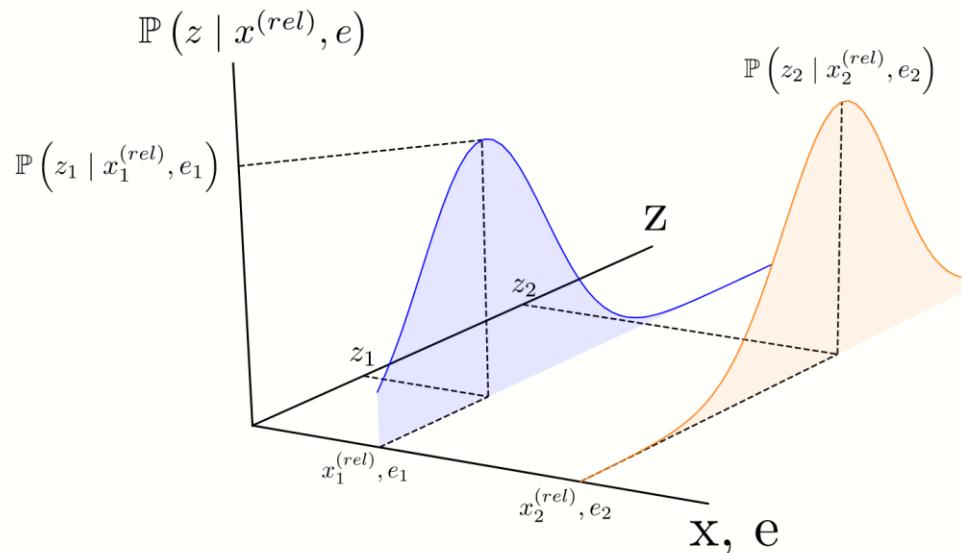
Example optimization path:



A Closer Look...

The learned model is constrained on a sparse set of triplets $\mathcal{Z}, \mathcal{X}^{(rel)}, \mathcal{E}$

⇒ Can try to improve inference by shaping the model elsewhere



Fitting the Viewpoint-Dependent Model – Take 2 (a)

Use a feature extractor to reduce factor dimensionality

$$\begin{aligned} \mathbb{P}(\mathcal{Z}_k \mid \mathcal{X}_k^{(rel)}, \mathcal{E}) &\doteq \mathbb{P}_{\theta}\left(f_{\psi}(\mathcal{Z}_k) \mid \mathcal{X}_k^{(rel)}, \mathcal{E}\right) \\ &\doteq \mathcal{N}\left(f_{\psi}(\mathcal{Z}_k); \mu_{\theta}(\mathcal{X}_k^{(rel)}, \mathcal{E}), \Sigma\right) \end{aligned}$$

In practice we use $\dim(\mathcal{E}) = 16$ $\dim(f_\psi(\mathcal{Z}_k)) = 12$.

Jointly fit all parameters

$$\arg \min_{\theta, \psi, \mathcal{E}_{1:n}} J(\theta, \psi, \mathcal{E}_{1:n}) = \arg \min_{\theta, \psi, \mathcal{E}_{1:n}} J_d(\theta, \psi, \mathcal{E}_{1:n}) + J_r(\mathcal{E}_{1:n}) + J_f(\psi)$$

Fitting the Viewpoint-Dependent Model – Take 2 (b)

Data term

$$J_d^{(i)}(\theta, \psi, \mathcal{E}_{1:n}) = \underset{\substack{\text{data term} \\ \text{for training example (i)}}}{\mathbb{E}}_{\Delta \mathcal{X}^{(rel)}, \Delta \mathcal{E}} \left\{ \frac{\| [\Delta \mathcal{X}^{(rel)}, \Delta \mathcal{E}] + s (\mathcal{X}^{(rel)} + \Delta \mathcal{X}^{(rel)}, \mathcal{E} + \Delta \mathcal{E}) \|^2}{\| [\Delta \mathcal{X}^{(rel)}, \Delta \mathcal{E}] \|^2} \right\}$$

where s is the optimization step at $\mathcal{X}^{(rel)} + \Delta \mathcal{X}^{(rel)}, \mathcal{E} + \Delta \mathcal{E}$ (that will be used at inference)

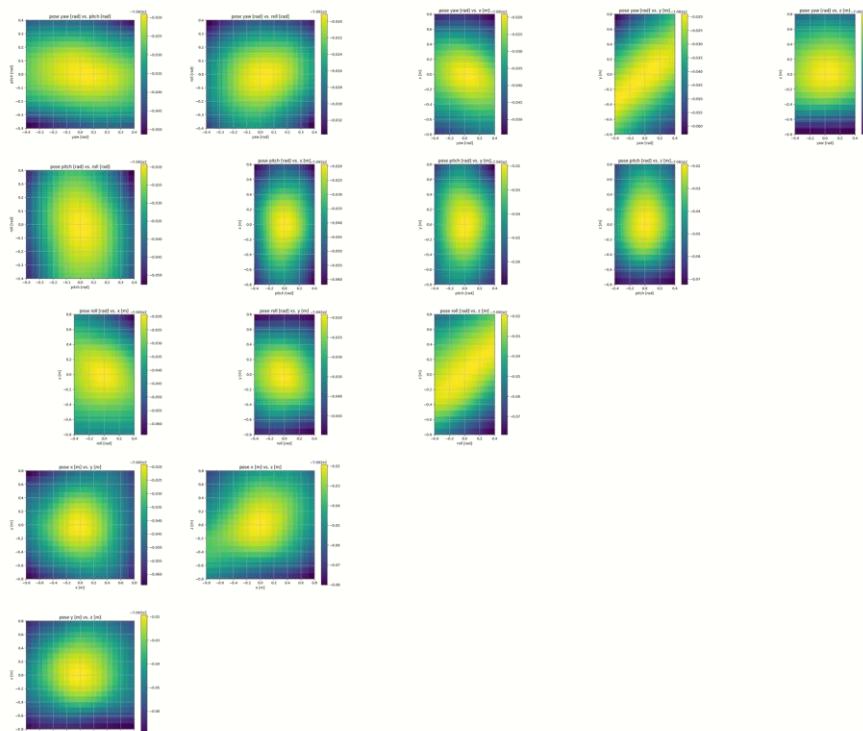
For gradient descent

$$s^{GD} \doteq -\frac{\partial}{\partial \mathcal{X}^{(rel)}, \mathcal{E}} \log \mathbb{P}_\theta (\mathcal{Z} \mid \mathcal{X}^{(rel)} + \Delta \mathcal{X}^{(rel)}, \mathcal{E} + \Delta \mathcal{E})$$

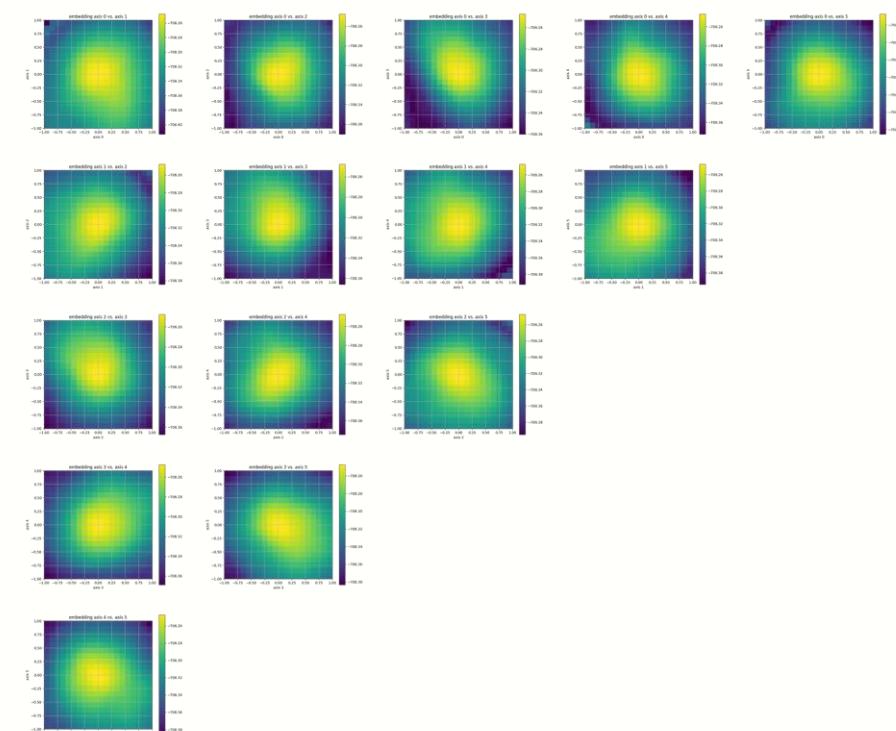
Learned Model

This finally gives useful gradients with respect to both $\mathcal{X}^{(rel)}$ and \mathcal{E}
Inference experiments still in progress.

likelihood - offsets in relative pose



likelihood - offsets in semantic representation



Outline

Viewpoint-Dependent models for Semantic Perception under Uncertainty

- ✓ 1. The semantic perception problem (Object-Level SLAM)
- ✓ 2. Viewpoint-dependent semantic measurement models
- ✓ 3. Contributions:
 - I. Classification under Model and Localization Uncertainty
 - II. Data Association-Aware Semantic Mapping and Localization
 - III. Semantic Perception with a Continuous Learned Representation
- 4. Summary

Summary

- ❖ We showed how to address semantic perception under uncertainty by exploiting the coupling between semantics and geometry provided by viewpoint-dependent models.

Contributions:

- I. Classification aware of Model uncertainty and correlations among viewpoints
- II. Data Association-Aware Semantic Mapping and Localization
- III. A novel approach to semantic SLAM through inference in a learned latent space

Thanks for listening!

Questions?