

A Generative Model for 3D Urban Scene Understanding from Movable Platforms

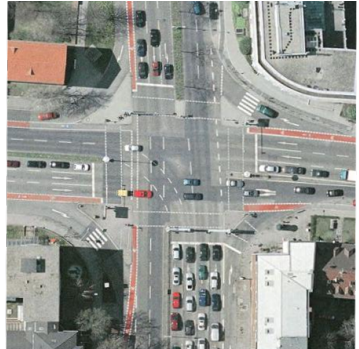
Andreas Geiger*, Martin Lauer*, Raquel Urtasun**

*KARLSRUHE INSTITUTE OF TECHNOLOGY

**TOYOTA TECHNOLOGICAL INSTITUTE AT CHICAGO



3D Urban Scene Understanding



Urban Scene Segmentation

- Semantic categories: *road, vehicles, building, sky*
- Joint detection and segmentation [Wojek 2008]
- Appearance combined with SfM [Sturgess 2009]

Geometric Methods

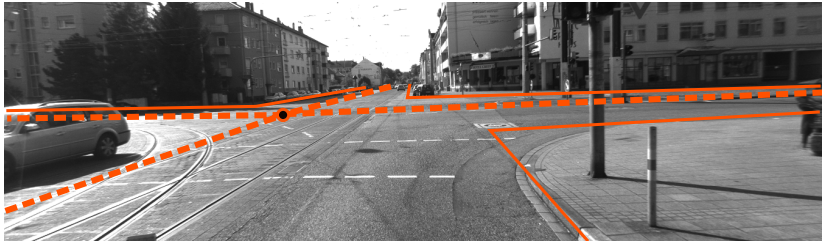
- 3D from single images [Hoiem 2007, Saxena 2009]
- Incorporating laws of physics [Gupta 2010]

Activity Recognition

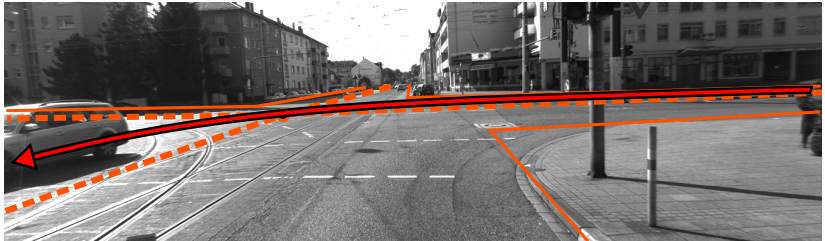
- Extraction of typical activity patterns in 2D optical flow
- Non-parametric clustering [Wang 2009, Kuettel 2010]



- **Goal:** Infer from short video sequences (moving observer)
 - **Topology and geometry** of the scene
 - **Semantic information** (e.g., traffic situation)
- **Probabilistic generative model** of 3D urban scenes
- **Static features:** Building facades
- **Dynamic features:** Moving vehicles



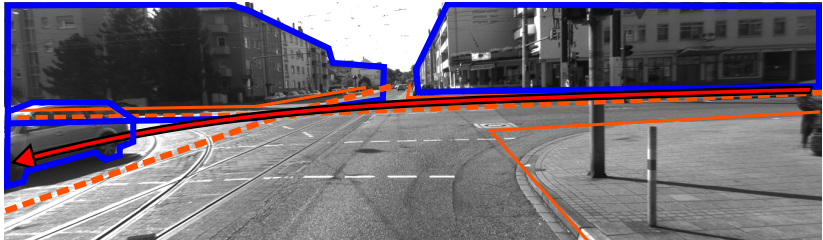
- **Goal:** Infer from short video sequences (moving observer)
 - **Topology** and **geometry** of the scene
 - **Semantic information** (e.g., traffic situation)
- Probabilistic generative model of 3D urban scenes
- Static features: Building facades
- Dynamic features: Moving vehicles



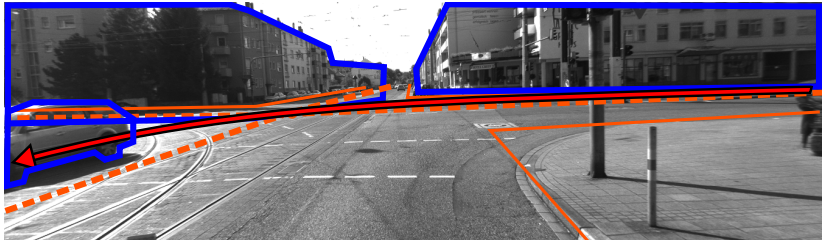
- **Goal:** Infer from short video sequences (moving observer)
 - **Topology** and **geometry** of the scene
 - **Semantic information** (e.g., traffic situation)
- Probabilistic generative model of 3D urban scenes
- Static features: Building facades
- Dynamic features: Moving vehicles



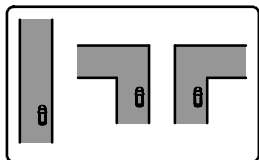
- **Goal:** Infer from short video sequences (moving observer)
 - **Topology** and **geometry** of the scene
 - **Semantic information** (e.g., traffic situation)
- **Probabilistic generative model** of 3D urban scenes
 - **Static features:** Building facades
 - **Dynamic features:** Moving vehicles



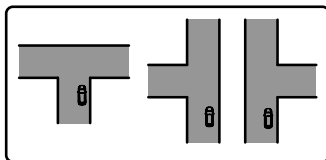
- **Goal:** Infer from short video sequences (moving observer)
 - **Topology** and **geometry** of the scene
 - **Semantic information** (e.g., traffic situation)
- **Probabilistic generative model** of 3D urban scenes
- **Static features:** Building facades
- **Dynamic features:** Moving vehicles



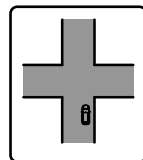
- **Goal:** Infer from short video sequences (moving observer)
 - **Topology** and **geometry** of the scene
 - **Semantic information** (e.g., traffic situation)
- **Probabilistic generative model** of 3D urban scenes
- **Static features:** Building facades
- **Dynamic features:** Moving vehicles



Topology A ($k = 2$)

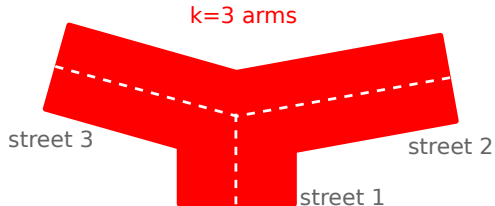


Topology B ($k = 3$)

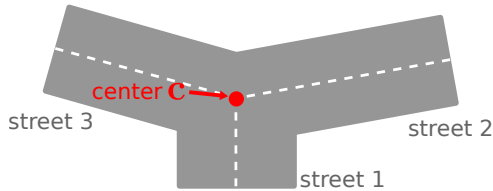


Topology C ($k = 4$)

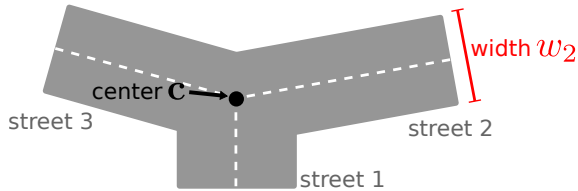
- Reasoning in **bird's eye perspective**



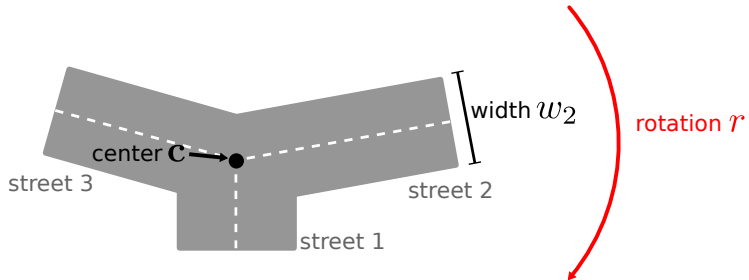
- k ... number of intersection arms
- \mathbf{c} ... center of intersection
- w ... street width
- r ... global rotation
- \mathbf{o} ... relative street orientation



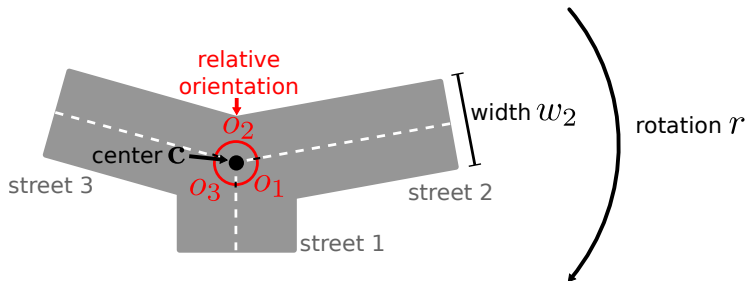
- k ... number of intersection arms
- \mathbf{c} ... center of intersection
- w ... street width
- r ... global rotation
- \mathbf{o} ... relative street orientation



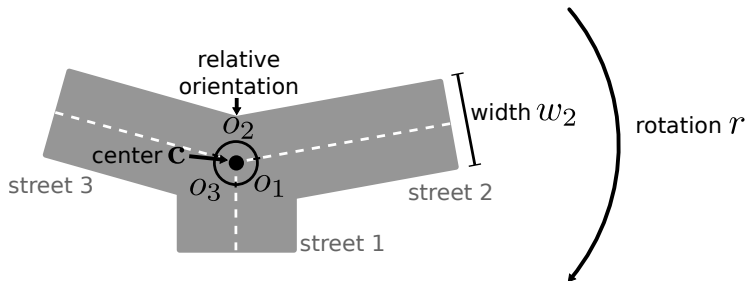
- k ... number of intersection arms
- \mathbf{c} ... center of intersection
- \mathbf{w} ... street width
- r ... global rotation
- \mathbf{o} ... relative street orientation



- k ... number of intersection arms
- \mathbf{c} ... center of intersection
- \mathbf{w} ... street width
- r ... global rotation
- \mathbf{o} ... relative street orientation



- k ... number of intersection arms
- \mathbf{c} ... center of intersection
- \mathbf{w} ... street width
- r ... global rotation
- \mathbf{o} ... relative street orientation

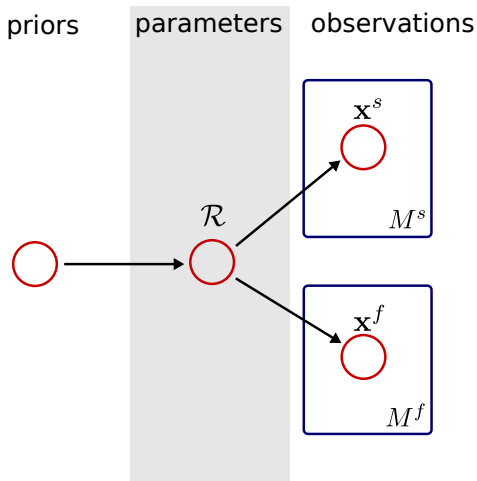


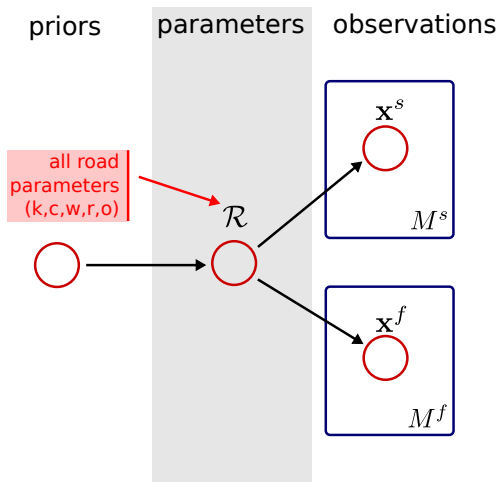
- k ... number of intersection arms
- \mathbf{c} ... center of intersection
- \mathbf{w} ... street width
- r ... global rotation
- \mathbf{o} ... relative street orientation

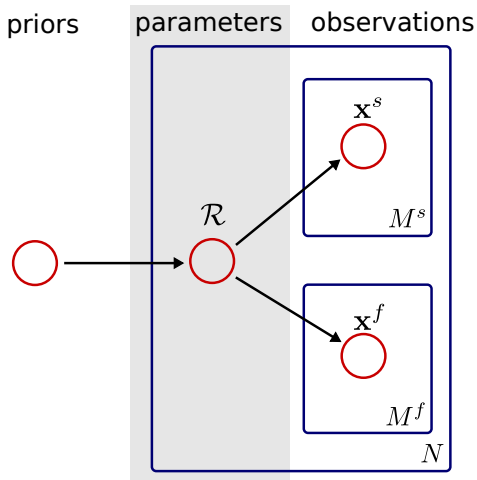
} \mathcal{R}

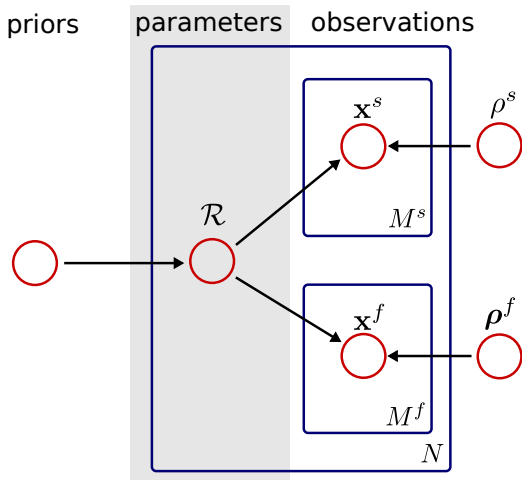


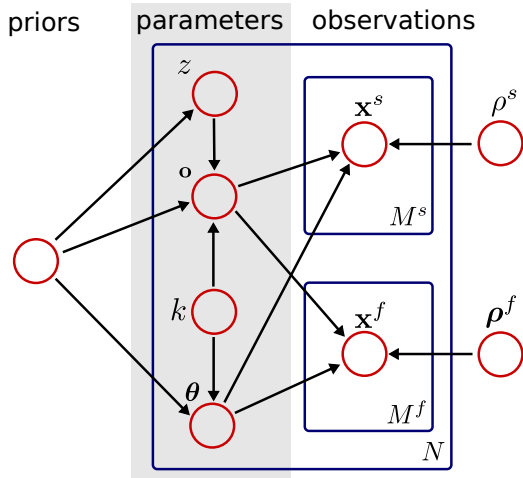
- **Static features:** Occupancy grids. For each cell, a variable indicates if it is occupied (+1), unobserved (0) or free (-1).
- **Dynamic features:** Sparse 3D scene flow
- Registration using **stereo visual odometry**

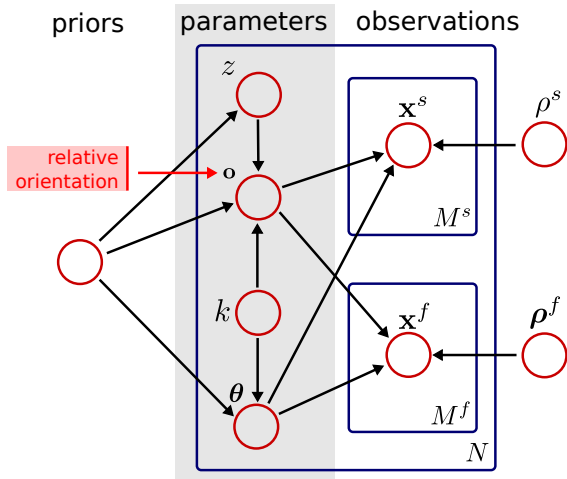


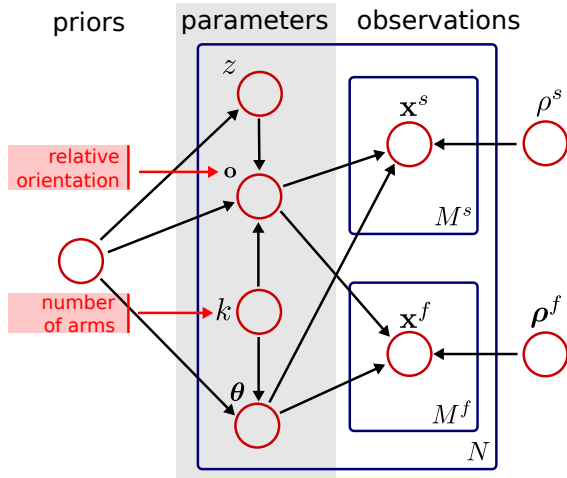


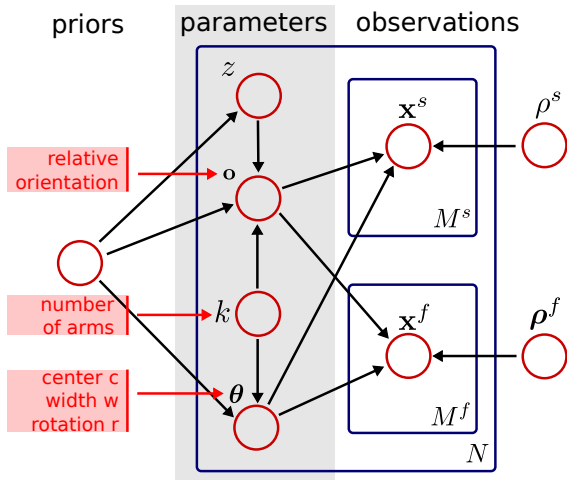




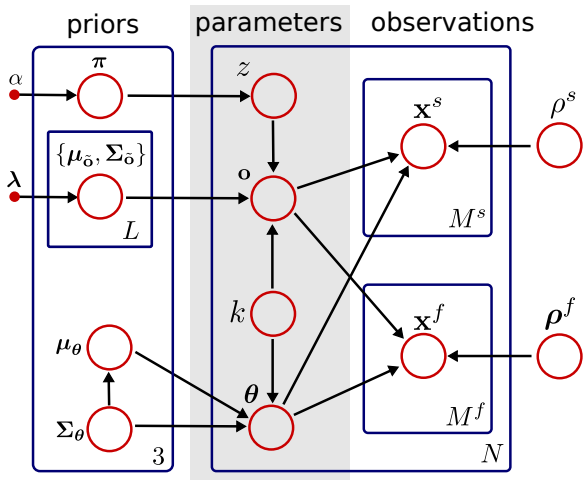




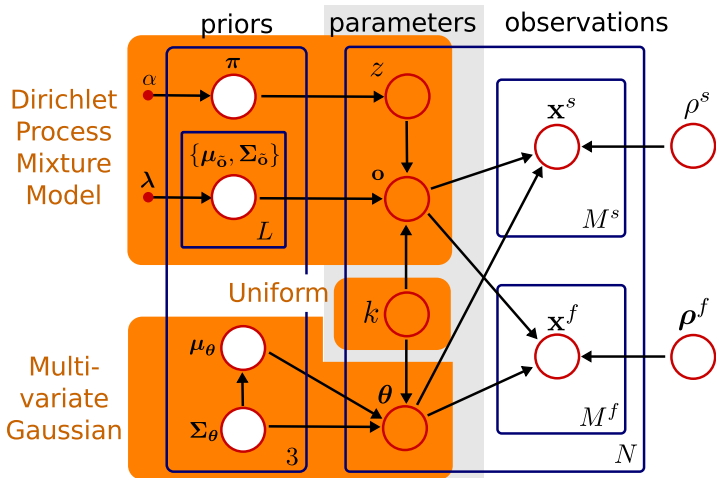


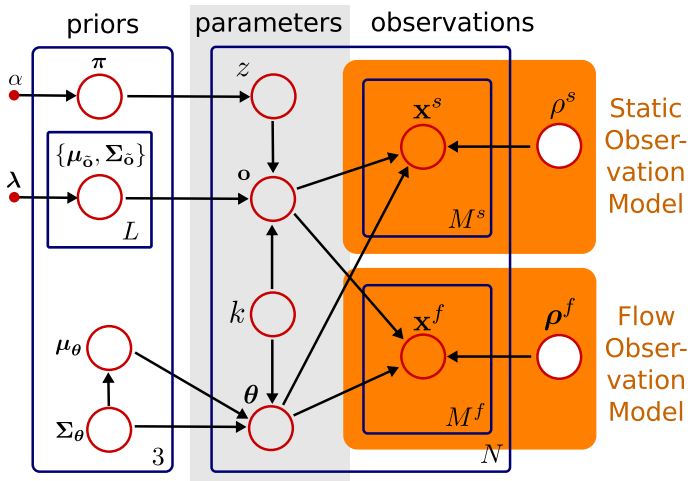


Probabilistic Model



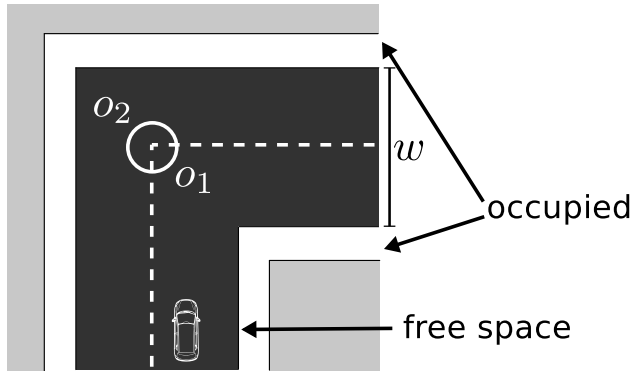
Probabilistic Model





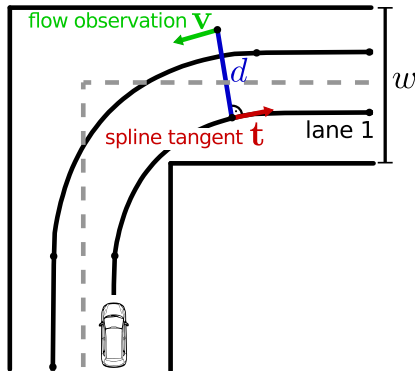
■ Static observation likelihood

$$p(\mathbf{x}^s | \mathbf{o}, \theta, \rho^s) \propto \exp\{\beta f(\mathbf{x}^s, \mathbf{o}, \theta, \rho^s)\}$$



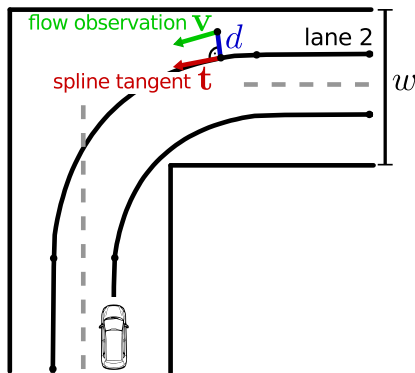
■ Flow observation likelihood

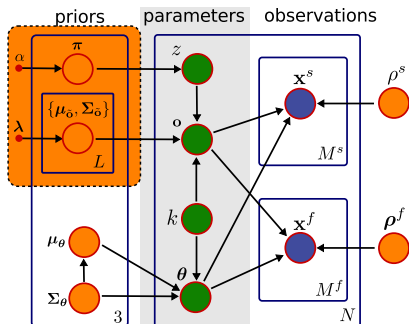
$$p(\mathbf{x}^f | \mathbf{o}, \theta, \rho^f) \propto \max_i \exp \left\{ -\frac{d_i^2}{2\rho_{f1}^2} - \frac{\|\mathbf{v} - \mathbf{t}_i\|_2^2}{2\rho_{f2}^2} \right\}$$



■ Flow observation likelihood

$$p(\mathbf{x}^f | \mathbf{o}, \theta, \rho^f) \propto \max_i \exp \left\{ -\frac{d_i^2}{2\rho_{f1}^2} - \frac{\|\mathbf{v} - \mathbf{t}_i\|_2^2}{2\rho_{f2}^2} \right\}$$





Learning:

■ Orientation:

Gibbs sampling with Dirichlet Process Mixture Model (MAP)

■ Center/rotation/width:

Maximum Likelihood

■ Observation model:

MH sampling (MAP)

Inference:

■ Reversible Jump MCMC:

■ Local MH moves

■ Global MH moves

■ Reversible jumps

Learning:

■ Orientation:

Gibbs sampling with Dirichlet Process Mixture Model (MAP)

■ Center/rotation/width:

Maximum Likelihood

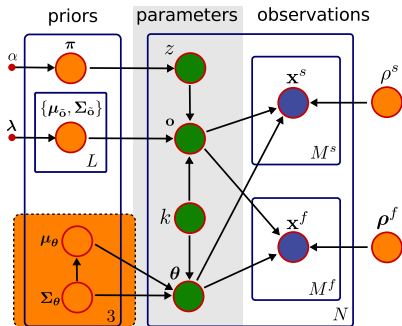
■ Observation model:

MH sampling (MAP)

Inference:

■ Reversible Jump MCMC:

- Local MH moves
- Global MH moves
- Reversible jumps



Learning:

■ Orientation:

Gibbs sampling with Dirichlet Process Mixture Model (MAP)

■ Center/rotation/width:

Maximum Likelihood

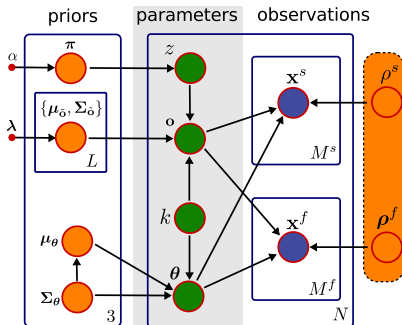
■ Observation model:

MH sampling (MAP)

Inference:

■ Reversible Jump MCMC:

- Local MH moves
- Global MH moves
- Reversible jumps



Learning:

■ Orientation:

Gibbs sampling with Dirichlet Process Mixture Model (MAP)

■ Center/rotation/width:

Maximum Likelihood

■ Observation model:

MH sampling (MAP)

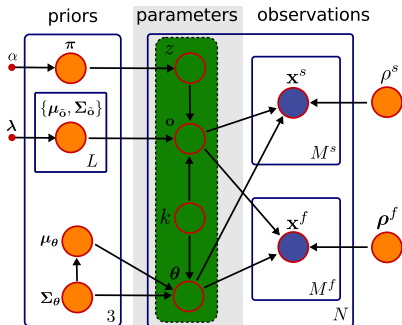
Inference:

■ Reversible Jump MCMC:

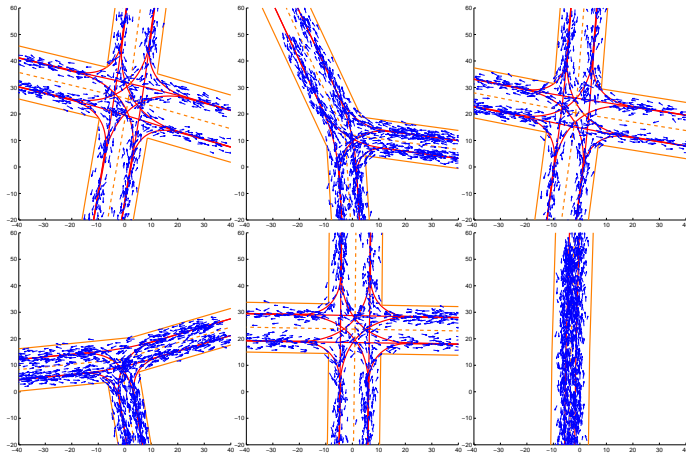
■ Local MH moves

■ Global MH moves

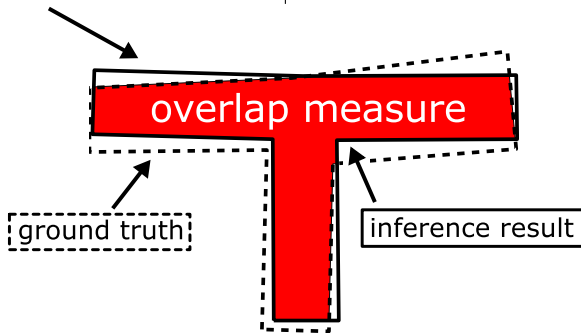
■ Reversible jumps



■ Samples from the prior:



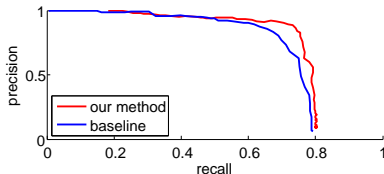
| | GP Regression | Ours |
|----------------------------|---------------|---------|
| Accuracy in estimating k | 61.1 % | 92.9 % |
| Location error | 5.4 m | 4.4 m |
| Orientation error | 14.1 deg | 6.6 deg |
| Overlap measure | 49.3 % | 62.7 % |

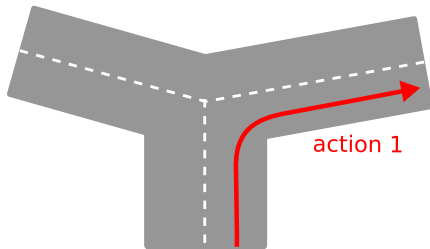


- **Re-weight scores** of [Felzenszwalb08] using spatial prior



- Increase in **average precision** from 71.3 % to 74.9 %

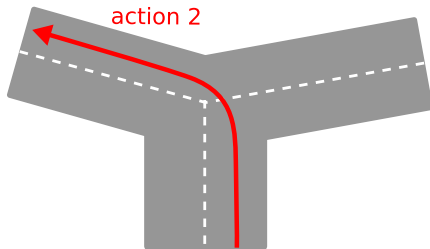




- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |

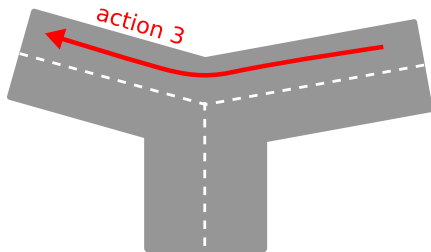
Extracting Semantic Activities



- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

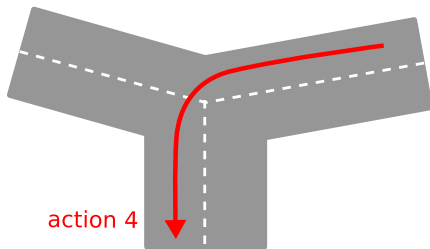
| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |

Extracting Semantic Activities



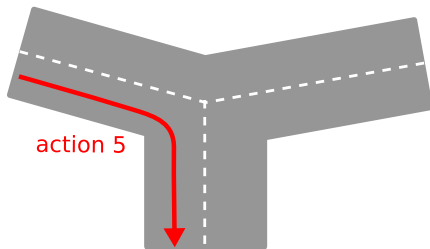
- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |



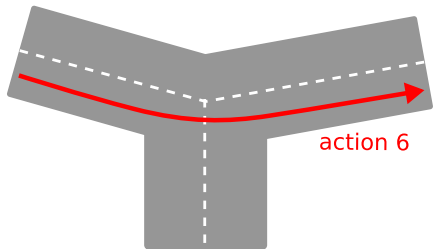
- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |



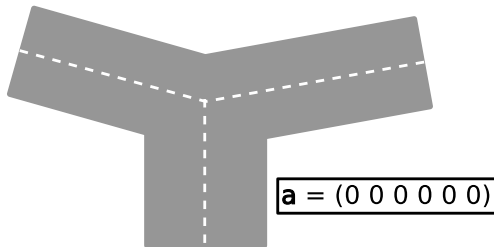
- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |



- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

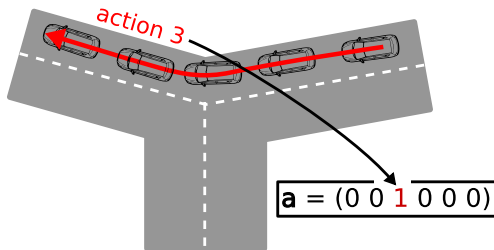
| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |



- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

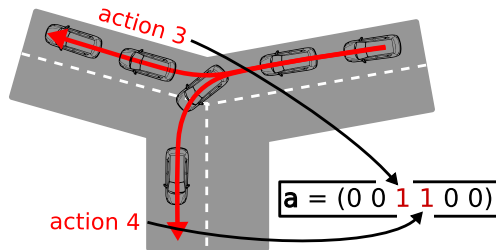
| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |

Extracting Semantic Activities



- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

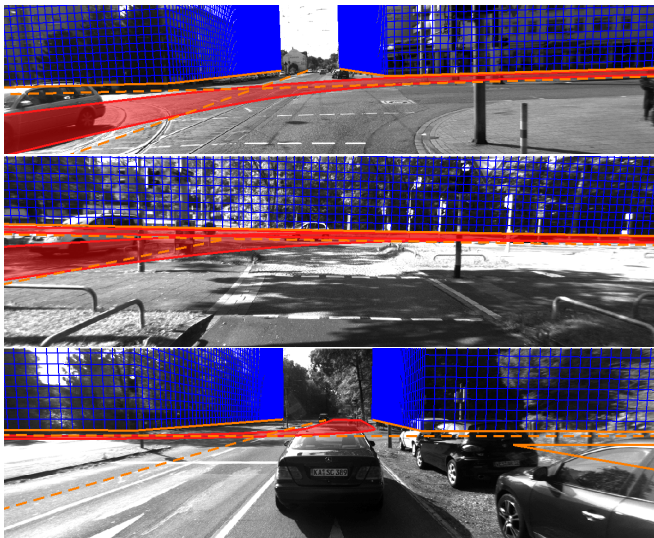
| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |



- **Activity:** $k(k - 1)$ dimensional binary vector \mathbf{a}
- Error measure: **Normalized Hamming distance**
- Results:

| | GP regression | Ours |
|------------------|---------------|------|
| Hamming distance | 0.16 | 0.08 |

Inference Results





Conclusion:

- **Generative model of 3D urban scenes**
- **Static** and **dynamic** features
- Improved **object detection & activity recognition**

Future work:

- **Features:** Vanishing points, scene labels, ...
- **Joint object and scene layout** reasoning
- Scene understanding with a **single camera**



Conclusion:

- **Generative model of 3D urban scenes**
- **Static** and **dynamic** features
- Improved **object detection & activity recognition**

Future work:

- **Features:** Vanishing points, scene labels, ...
- **Joint object and scene layout** reasoning
- Scene understanding with a **single camera**



Conclusion:

- **Generative model of 3D urban scenes**
- **Static** and **dynamic** features
- Improved **object detection & activity recognition**

Future work:

- **Features:** Vanishing points, scene labels, ...
- **Joint object and scene layout** reasoning
- Scene understanding with a **single camera**



Conclusion:

- **Generative model** of **3D urban scenes**
- **Static** and **dynamic** features
- Improved **object detection & activity recognition**

Future work:

- **Features:** Vanishing points, scene labels, ...
- **Joint object and scene layout** reasoning
- Scene understanding with a **single camera**



Conclusion:

- **Generative model of 3D urban scenes**
- **Static** and **dynamic** features
- Improved **object detection & activity recognition**

Future work:

- **Features:** Vanishing points, scene labels, ...
- **Joint object and scene layout** reasoning
- Scene understanding with a **single camera**



Conclusion:

- **Generative model** of **3D urban scenes**
- **Static** and **dynamic** features
- Improved **object detection** & **activity recognition**

Future work:

- **Features:** Vanishing points, scene labels, ...
- **Joint object and scene layout** reasoning
- Scene understanding with a **single camera**

Unique orientation vector \mathbf{o} , constrained to the Δ^{k-1} simplex

$$\sum_{i=1}^k o_i = 1 \quad o_i \geq 0 \quad \alpha_i = r + 2\pi \sum_{j=1}^{i-1} o_j$$

