

# A Generative Model for Online Depth Fusion

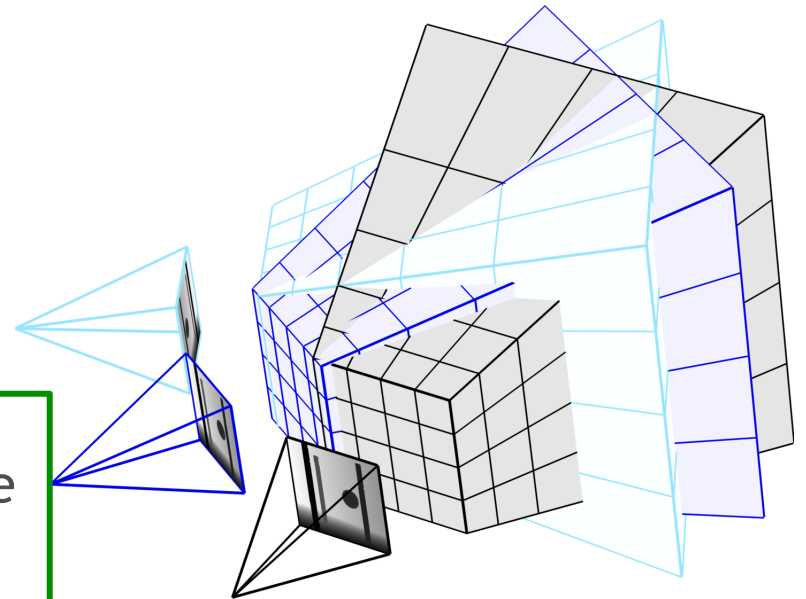
George Vogiatzis (Aston University) and  
Oliver Woodford (Toshiba Research)

# Background

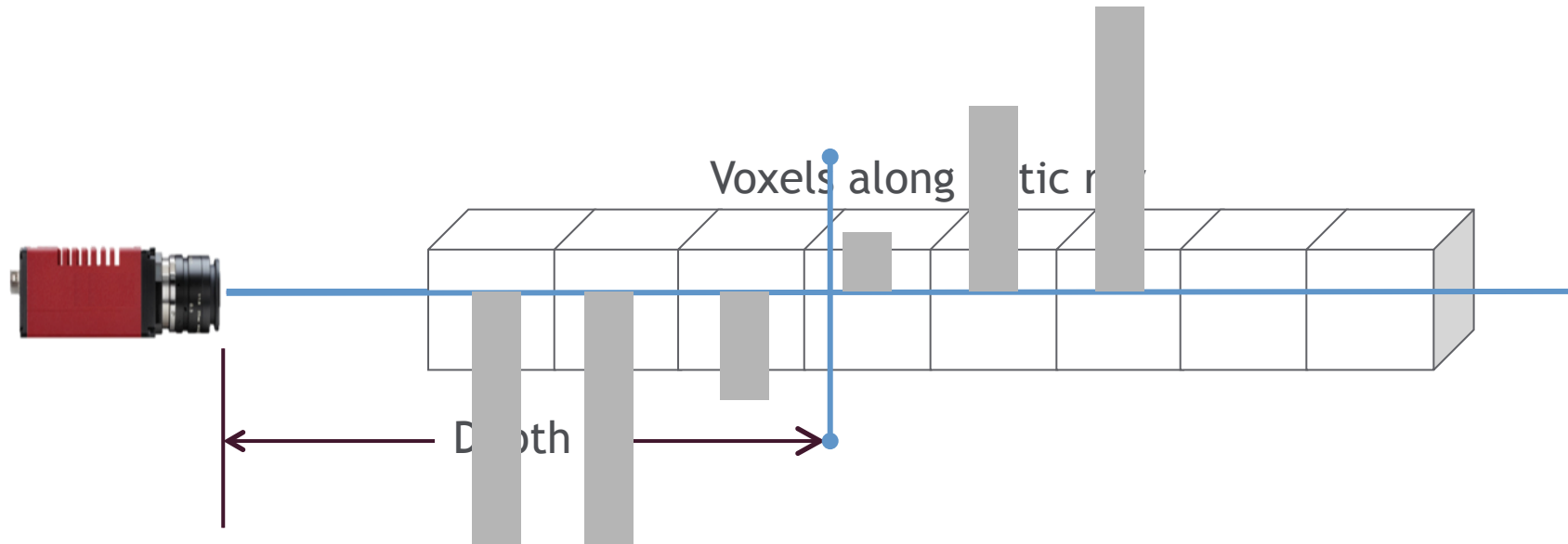
- ▶ Plethora of depth-measuring technologies
  - ▶ binocular/multi-view stereo
  - ▶ Structured light stereo (e.g. Kinect)
  - ▶ Sonar
  - ▶ Time-of-flight
  - ▶ Laser
  - ▶ ...

# Depth-map fusion

- ▶ Convert depth-maps to scene geometry
  - ▶ Crucial problem
- ▶ Offline fusion: collect all depth-maps THEN merge
  - ▶ Point-cloud
  - ▶ Octree
  - ▶ TSDF
- ▶ Online fusion: merge each incoming depth-map into state
  - ▶ Forward/inverse sensor modeling (Robotics)
  - ▶ TSDF (KinectFusion)



# Truncated Signed Distance Functions



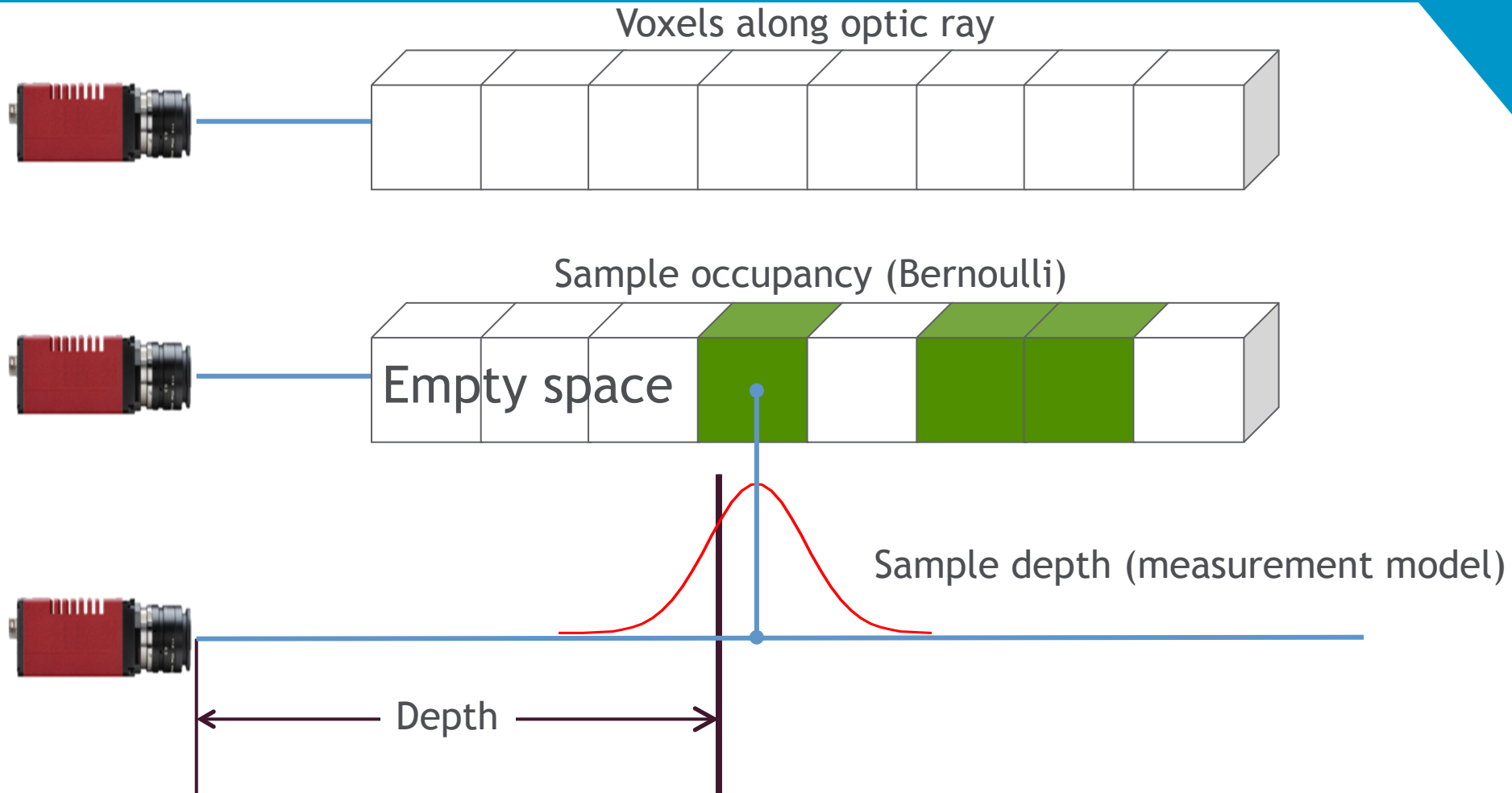
- ▶ Shown to be equivalent to accumulating probabilistic *evidence* of visibility (log-odds)
- ▶ Under a logistic sensor noise model ( $\text{sech}(x)^2$ )
- ▶ No account of outlier measurements



# Robotics online depth fusion

- ▶ Inverse models (Elfes & Matthies 87, Konolige 97)
  - ▶ Model directly  $p(\text{occupancy} | \text{depth})$
  - ▶ No inter-dependency of occupancy variables along an optic ray (free-space constraints)
  
- ▶ Forward models (Thrun 01, Pathak 07)
  - ▶ Model  $p(\text{depth} | \text{visibility})$
  - ▶ Assume occupancy is driven by visibility
  - ▶ Cannot model occlusion

# Generative model of depth measurement

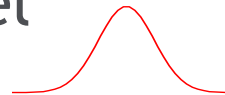


# Generative model of depth measurement

- ▶ Occupancy  $\mathbf{x} = \{x_i\}_{i=1}^N, x_i \in \{0, 1\}$
- ▶ Depth measurement  $y \in \mathbb{R}$
- ▶ Visible voxel index  $v = 1, 2, \dots$

$$p(y|\mathbf{x}) = \sum_v p(y|v) p(v|\mathbf{x})$$

Measurement  
model



Is 1 iff  $v$  is index  
of first occupied  
voxel

# Inferring occupancy

- ▶ Prior  $p(x)$  factorizes
- ▶ BUT, posterior  $p(x|y)$  doesn't!
- ▶ Not feasible to maintain full covariance in online fusion
- ▶ Our approach:
  - ▶ Assume factored approximation  $q(x)$
  - ▶ Minimize  $KL(p||q)$  (Expectation Propagation)
  - ▶ Amounts to computing  $p(x_i|y)$  marginals
    - ▶ **Tractable!**

# Inferring occupancy

- ▶ Depth likelihood depends on occupancies through index of visible voxel

$$p(x_i|y) = \sum_v p(x_i v|y)$$

$$= \sum_v \underbrace{p(x_i|v)} \underbrace{p(v|y)}$$

$$= \underbrace{p(v = i|y)} + p(x_i) \sum_{v=1}^{i-1} \underbrace{p(v|y)}$$

$$p(v) = p(x_v) \prod_{i=1}^{v-1} (1 - p(x_i))$$

$$p(v|y) = p(y|v) \frac{p(v)}{p(y)}$$

$$p(x_i|v) = \begin{cases} 1 & v = i \\ 0 & v > i \\ p(x_i) & v < i \end{cases}$$

# Outlier measurements

- ▶ Can use a simple noise+outlier model

$$p(y|v)$$

# Outlier measurements

- ▶ Can use a simple noise+outlier model

$$p(y|v, \omega) = \omega \cdot \mathcal{C}(y) + (1 - \omega) \cdot \mathcal{M}_v(y)$$

Outlier ratio      Clutter dist.      Noise dist.

- ▶ Assume  $\omega$  comes from Beta( $\alpha, \beta$ ) hyper-prior

# Full model

- ▶ Factorized prior  $p(\mathbf{x}, \omega) = p(\mathbf{x}) p(\omega)$

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i | \gamma_i), \quad p(x | \gamma) = \gamma^x (1 - \gamma)^{1-x},$$

$$p(\omega) = \prod_{i=1}^N p(\omega_i | \alpha_i, \beta_i), \quad p(\omega | \alpha, \beta) = \frac{\omega^{\alpha-1} (1 - \omega)^{\beta-1}}{B(\alpha, \beta)}$$

- ▶ Noise + outlier likelihood

$$p(y | v, \omega) = \omega \cdot \mathcal{C}(y) + (1 - \omega) \cdot \mathcal{M}_v(y)$$



# Inference

- ▶ Assume factored approximation

$$q(\mathbf{x}, \boldsymbol{\omega}) = \prod_{i=1}^N q_i(x_i) \prod_{j=1}^N q_j(\omega_j)$$

- ▶ Minimize KL divergence between

$$p(\mathbf{x}, \boldsymbol{\omega} | y) \text{ and } q(\mathbf{x}, \boldsymbol{\omega})$$

- ▶ Matching sufficient statistics (~EP)

$$\mathbb{E}_{q(\mathbf{x}, \boldsymbol{\omega})} [x_i] = \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\omega} | y)} [x_i],$$

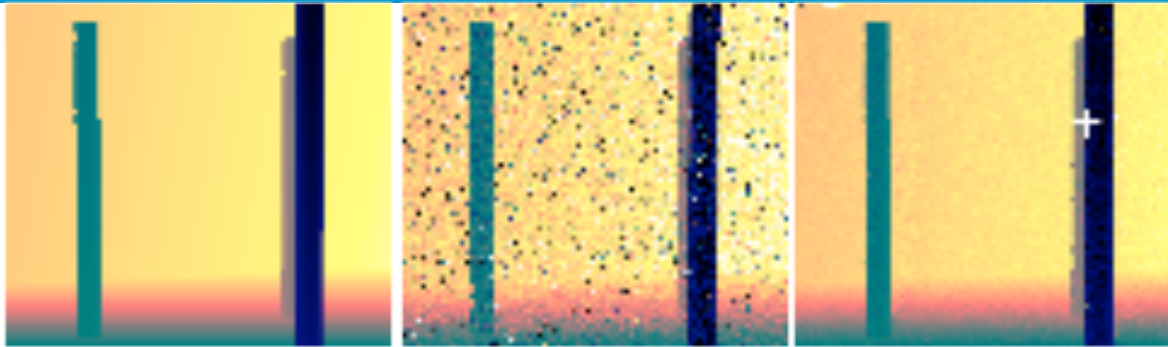
$$\mathbb{E}_{q(\mathbf{x}, \boldsymbol{\omega})} [\ln \omega_i] = \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\omega} | y)} [\ln \omega_i],$$

$$\mathbb{E}_{q(\mathbf{x}, \boldsymbol{\omega})} [\ln(1 - \omega_i)] = \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\omega} | y)} [\ln(1 - \omega_i)]$$

# Modelling options

- ▶ Outlier ratio
  - ▶ One fixed  $\omega$  (generative1)
  - ▶ One  $\omega$  that is estimated from data (generative2)
  - ▶ Multiple  $\omega$ , one per optic ray or per voxel, estimated from data (generative3)
- ▶ Online fusion
  - ▶ First order independence assumptions plus appearing and disappearing surfaces
  - ▶ Similar to ‘forgetting factor’ in TSDF
- ▶ How to output depth: use  $q(v)$

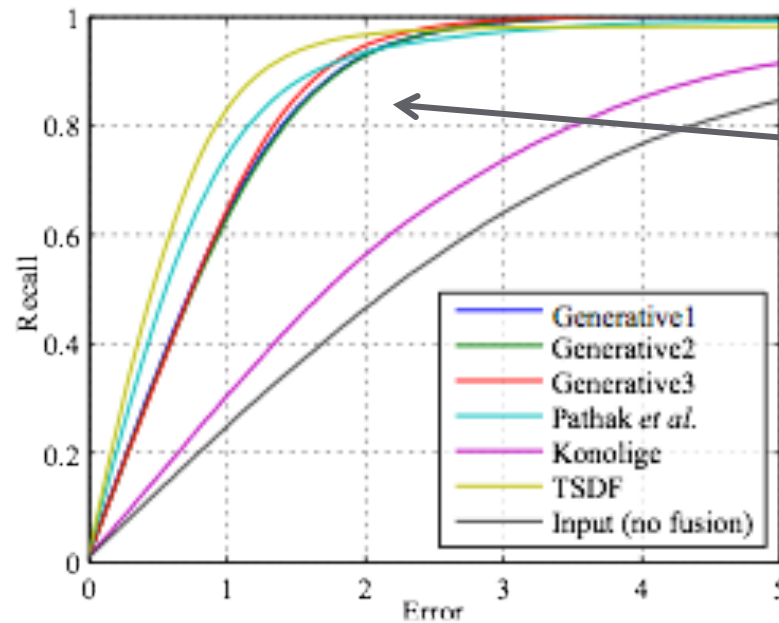
# Ground truth benchmarking



Ground truth

Noise added

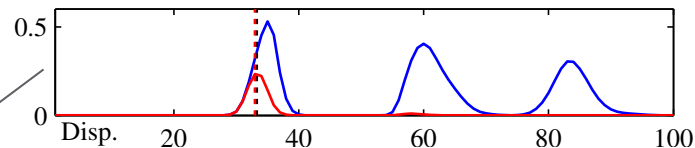
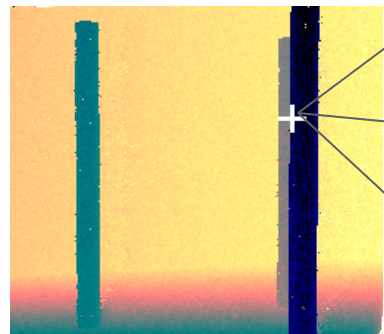
Fused output



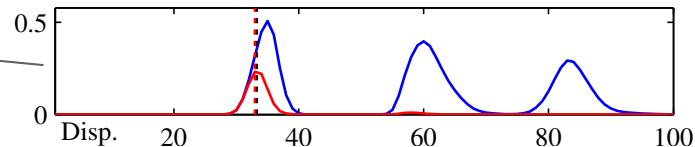
Score=area under curve normalized to 1.0

Error/recall curve

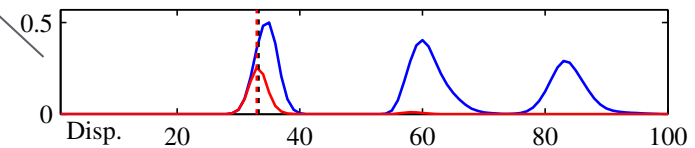
# Occupancy vs visible surface along a ray



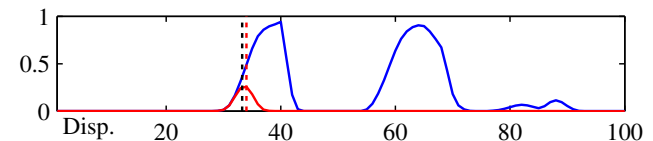
(a) Generative1



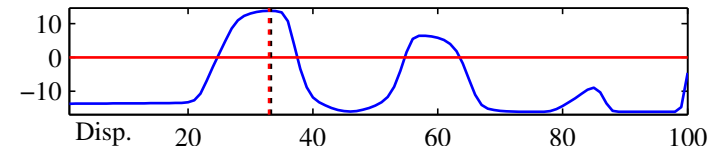
(b) Generative2



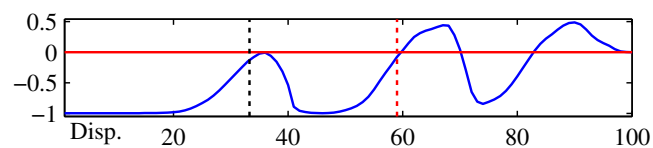
(c) Generative3



(d) Pathak *et al.*



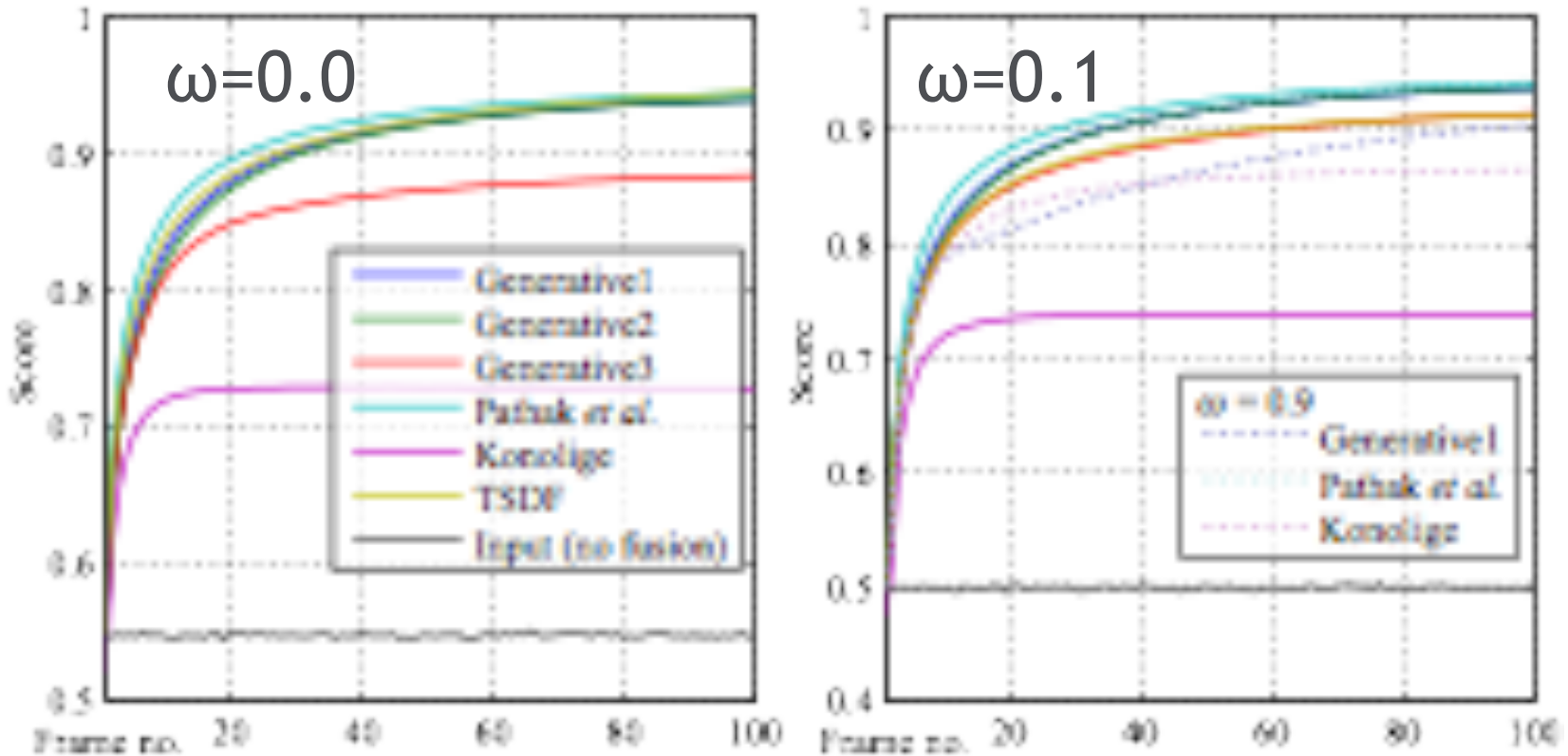
(e) Konolige



(f) TSDF

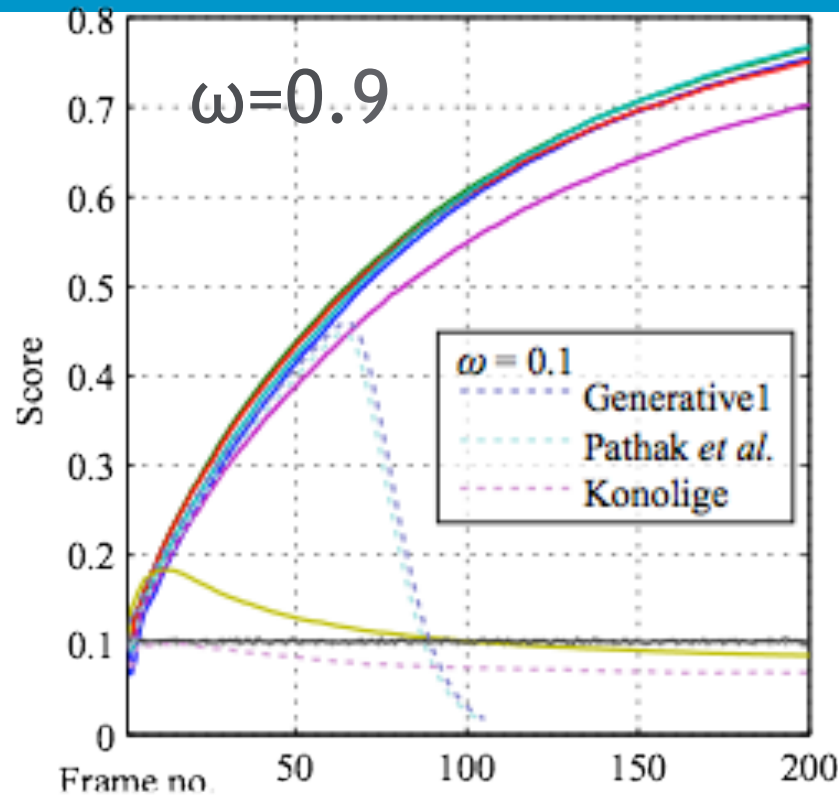
- ▶ More meaningful occupancy values, with better defined maxima

# Effect of outlier modelling



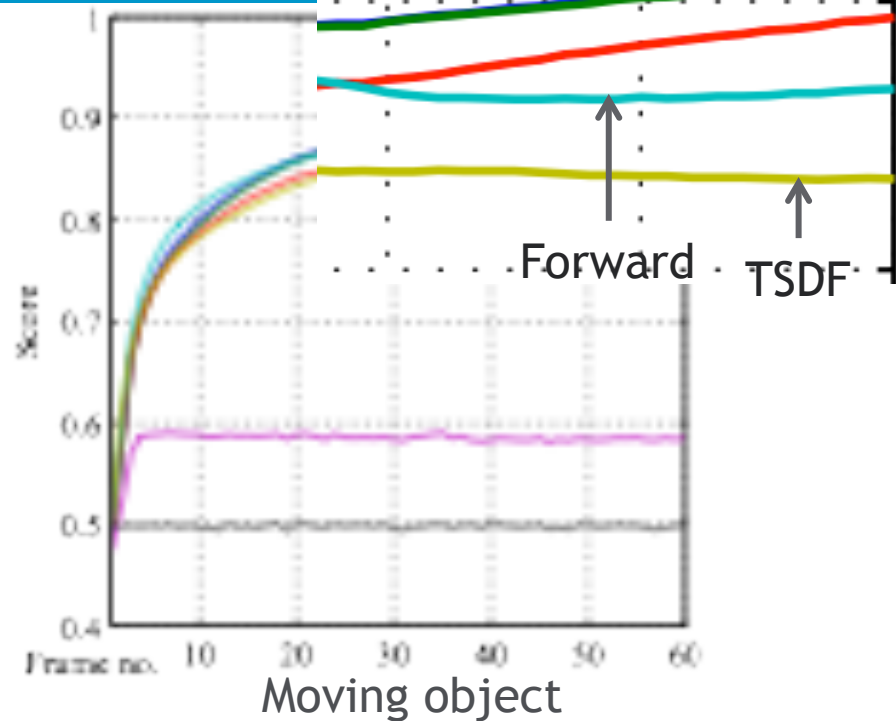
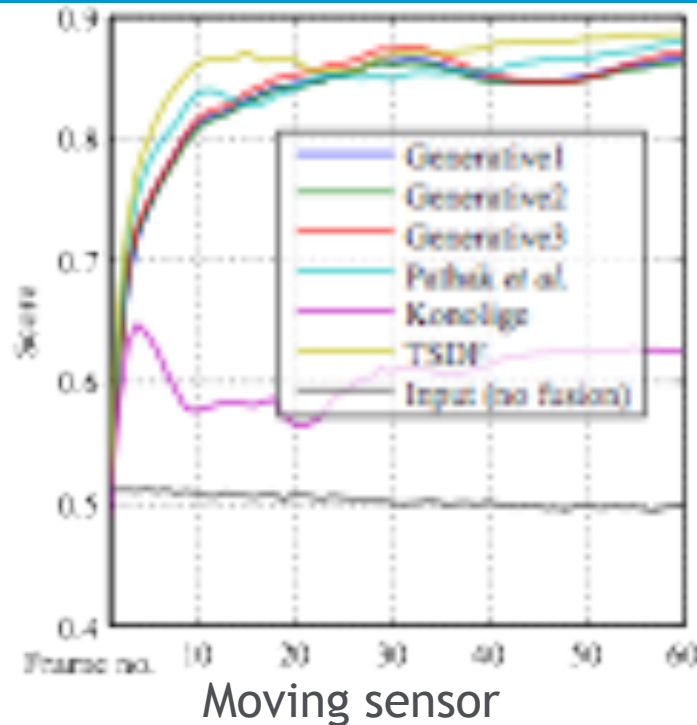
- ▶ Not much effect in low outlier regimes
- ▶ Generative2 (estimate one  $\omega$ ) is good performance/ computation compromise

# Effect of outlier modelling



- ▶ Quite significant in heavy-outlier regime
- ▶ TSDF cannot cope
- ▶ Methods that **model** outliers but do not **estimate**  $\omega$  also do poorly

# Motion



- ▶ Faster reaction time to occlusions/disocclusions
- ▶ Due to more accurate occlusion reasoning

# Ground truth benchmarking

Our method

Fixed  $\omega$

Estimate one  $\omega$

Estimate multiple  $\omega$

Forward model

Inverse model

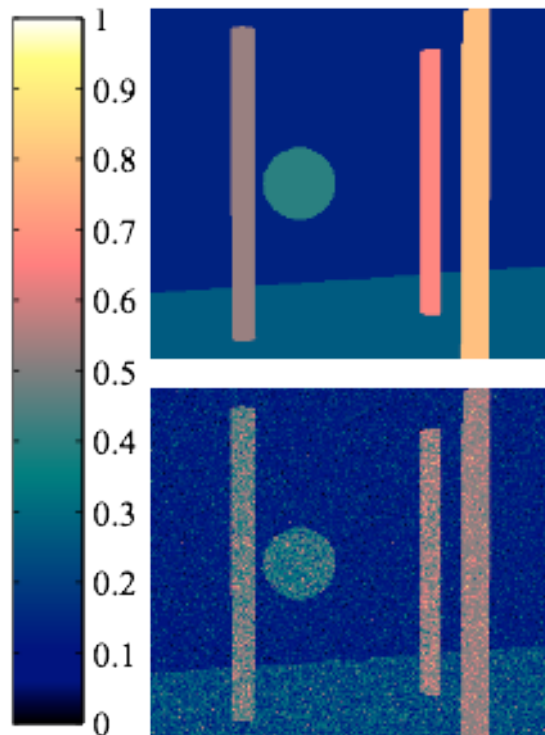
Kinect fusion

	$\omega$	Generative1	Generative2	Generative3	Pathak <i>et al.</i>	Konolige	TSDf
Static sensor	0	0.065±0.34	0.069±0.33	0.36±0.62	<b>0.000</b> ±0.33	1.37±0.25	-0.17±0.23
	0.1	0.064±0.35	0.059±0.34	0.19±0.45	<b>-0.002</b> ±0.33	1.32±0.25	-0.37±0.31
	0.4	0.090±0.39	0.073±0.39	0.070±0.40	<b>-0.002</b> ±0.37	1.18±0.33	-7.33±18.0
	0.9	0.60±7.08	<b>0.24</b> ±6.79	0.45±6.94	0.79±7.91	3.15±11.0	-23.7±20.0
Moving sensor	0	0.35±0.97	0.33±1.09	0.35±1.43	<b>0.27</b> ±0.86	2.09±8.49	-0.59±3.53
	0.1	0.35±2.13	<b>0.34</b> ±2.16	0.36±2.20	0.42±2.77	3.69±14.1	-0.69±4.35
	0.4	0.50±3.85	<b>0.41</b> ±4.09	0.48±3.87	0.74±4.86	6.90±21.1	-1.52±6.87
	0.9	<b>2.96</b> ±11.1	8.00±18.4	3.39±12.7	3.11±11.1	20.9±32.2	-10.1±12.6

- ▶ Smaller errors at high  $\omega$  regime
- ▶ TSDf struggles



# Estimating different outlier ratios



Ground truth  $\omega$

Estimated  $\omega$

- ▶ Can identify regions in the scene that produce different sensor response (e.g. shiny/textureless)
- ▶ First stage of scene classification scheme

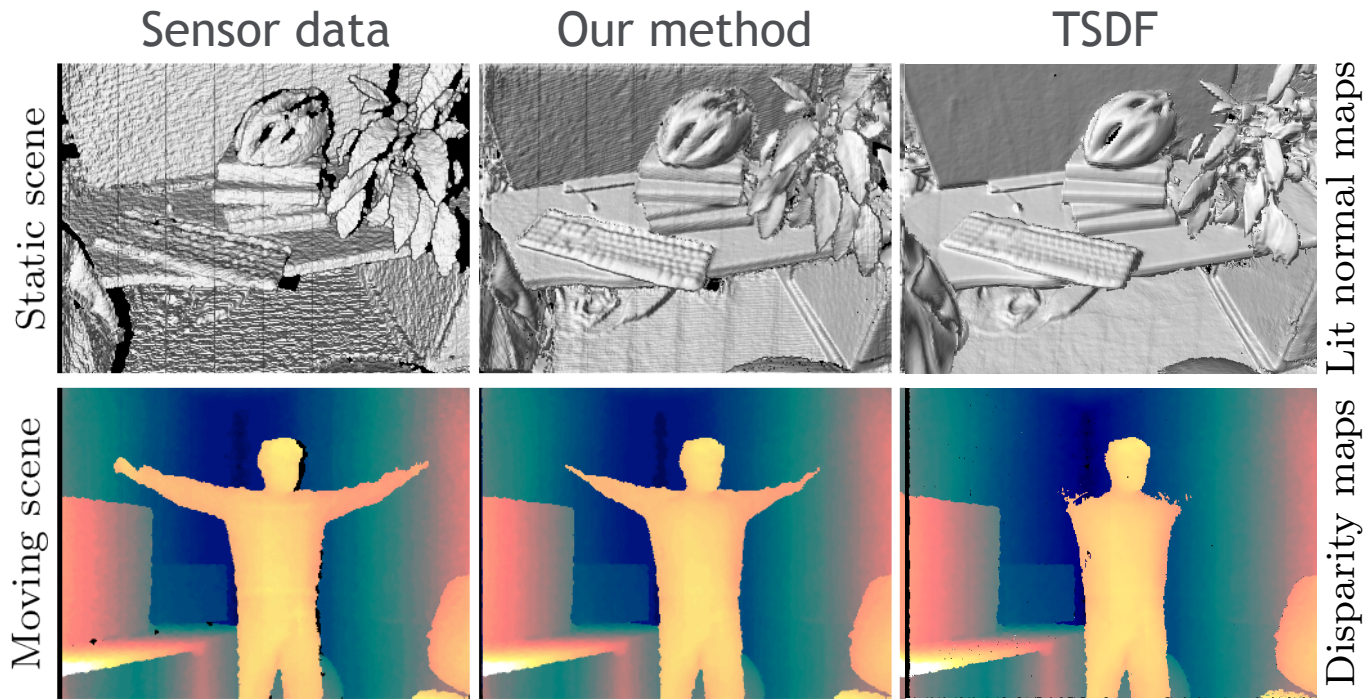
# Kinect - static scene/moving sensor



# Kinect - moving scene/static sensor



# Kinect data



# Multi-view stereo data

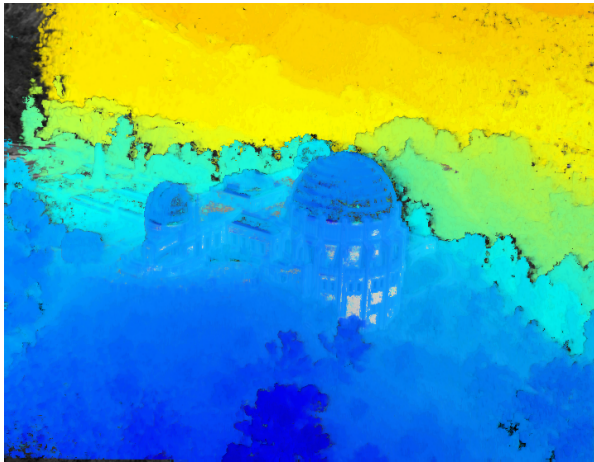
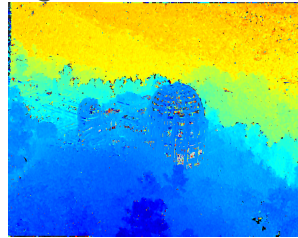


# Multi-view stereo data

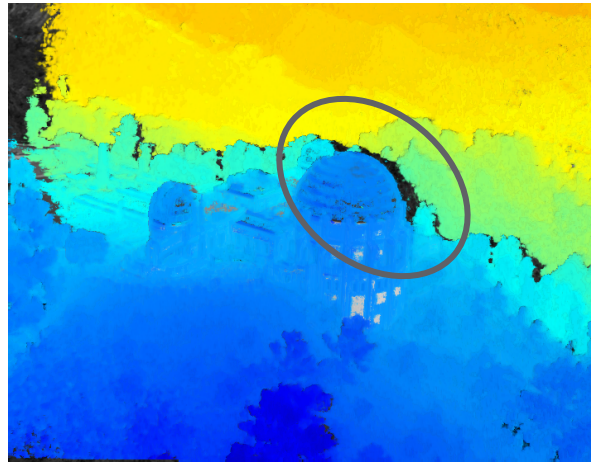
Input



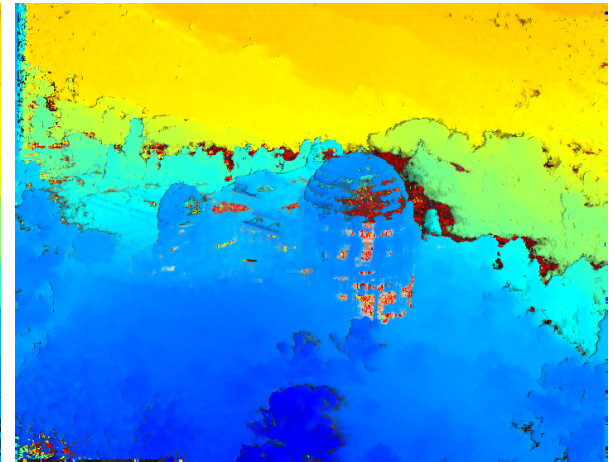
Estim. Depth using  
Vogiatzis&Hernandez 2011



Our method



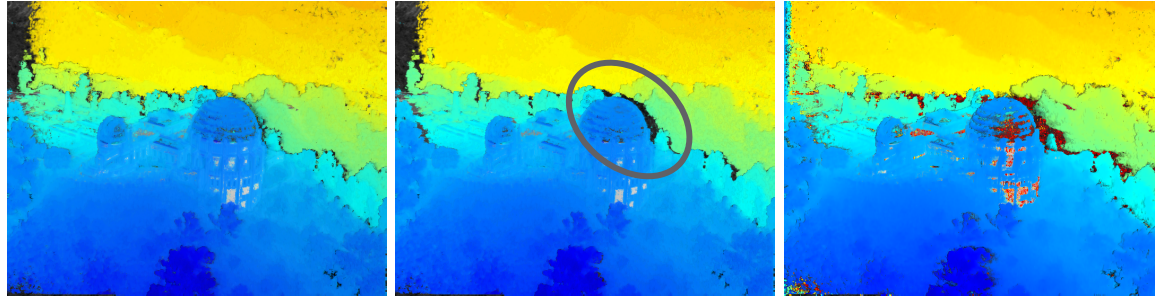
Forward model



TSDF



# Take home messages



Our method

Forward model

TSDF

- ▶ Use more realistic sensor modelling (e.g. outliers)
  - ▶ Generative models react faster to scene changes
  - ▶ Our occlusion model => better reconstruction at depth discontinuities
  - ▶ Inferring outlier ratios helps, but significantly only in extreme cases
    - ▶ interesting potential for scene classification
  - ▶ Volume resampling is slow (~9fps in KinectFusion) but allows zooming
- Many more experimental results in ECCV'12 paper
- Thanks!