# Joint Multiview Segmentation and Localization of RGB-D Images using Depth-Induced Silhouette Consistency

Chi Zhang[1], Zhiwei Li[2], Rui Cai[2], Hongyang Chao[1], Yong Rui[2]

1 Sun Yat-Sen University, Guangzhou, P.R. China

2 Microsoft Research, Beijing, P.R. China

Microsoft Research

## Motivation

- *Aim at RGB-D SLAM for object scanning.*
  - *Input: An RGB-D stream.*
  - *Output: foreground masks and camera poses of the extracted keyframes.*
- *Use silhouettes to improve pose estimation.*
- *Silhouettes are difficult to obtain in practice. Make it practical by generating silhouettes on-the-fly.*
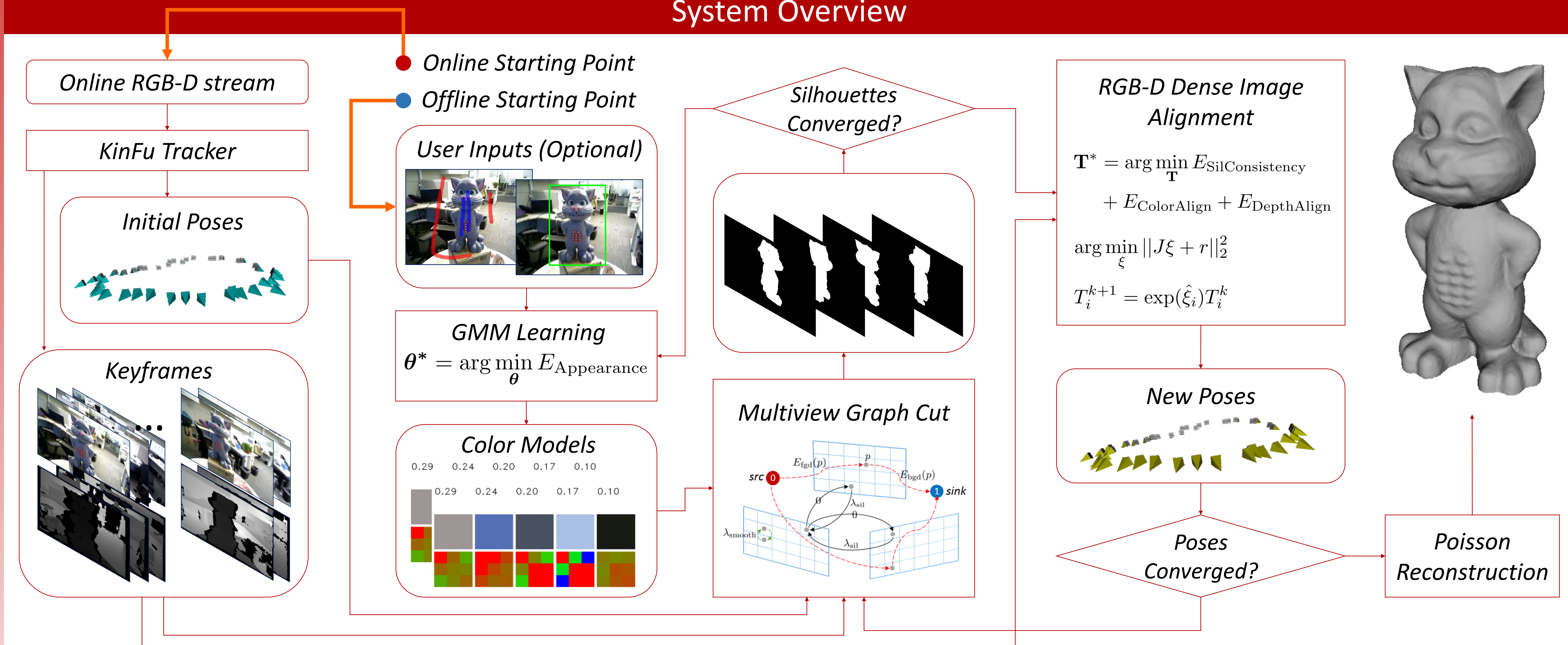
## Variables

*For each view i, we optimize:*

$\mathbf{S} \equiv \{S_i\}_{i=1}^N$    *A binary foreground mask.*

$\boldsymbol{\theta} \equiv \{\theta^{\text{fgd}}, \theta^{\text{bgd}}\}_{i=1}^N$    *Two GMM color models.*

$\mathbf{T} \equiv \{T_i\}_{i=1}^N$    *A local-to-world camera pose.*

## System Overview



Online RGB-D stream

KinFu Tracker

Initial Poses

Keyframes

- Online Starting Point
- Offline Starting Point

User Inputs (Optional)

GMM Learning
$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} E_{\text{Appearance}}$

Color Models

0.29 0.24 0.20 0.17 0.10
0.29 0.24 0.20 0.17 0.10

Multiview Graph Cut

Silhouettes Converged?

RGB-D Dense Image Alignment

$\mathbf{T}^* = \arg\min_{\mathbf{T}} E_{\text{SilConsistency}} + E_{\text{ColorAlign}} + E_{\text{DepthAlign}}$

$\arg\min_{\xi} ||J\xi + r||_2^2$

$T_i^{k+1} = \exp(\hat{\xi}_i) T_i^k$

New Poses

Poses Converged?

Poisson Reconstruction

## Formulation

*Five Sub-energies:*

$E_{\text{All}}(\mathbf{S}, \boldsymbol{\theta}, \mathbf{T}) = E_{\text{Appearence}}(\mathbf{S}, \boldsymbol{\theta})$

$+ E_{\text{MaskSmooth}}(\mathbf{S})$

$+ E_{\text{SilConsitency}}(\mathbf{S}, \mathbf{T})$

$+ E_{\text{ColorAlign}}(\mathbf{T})$

$+ E_{\text{DepthAlign}}(\mathbf{T})$

*Detailed definitions:*

$\longrightarrow \sum_i \sum_{p \in \Omega_i} -\text{Prob}(I_i(p) \mid S_i(p), \theta_i^{\text{bgd}}, \theta_i^{\text{fgd}})$

$\longrightarrow \sum_i \sum_{p,r \in \mathcal{N}_4} w_{pr} ||S_i(p) - S_i(r)||^2$

$\longrightarrow \sum_i \sum_{p \in \tilde{\Omega}_i} \sum_{j \neq i} S_i(p) \cdot ||S_i(p) - S_j(q)||^2$

$\longrightarrow \sum_i \sum_{p \in \tilde{\Omega}_i} \sum_{j \in \mathcal{N}_i} ||I_i(p) - I_j(q)||^2$

$\longrightarrow \sum_i \sum_{p \in \tilde{\Omega}_i} \sum_{j \in \mathcal{N}_i} ||D_i(p) - D_j(q)||^2$

*Optimize w.r.t.* $\mathbf{T}$ *by Gauss-Newton Method.*

*Optimize w.r.t.* $\mathbf{S}, \theta$ *by Multiview Graph Cut.*

$q = \pi_j(T_j^{-1} T_i \pi_i^{-1}(p, D_i(p)))$

$\tilde{\Omega}_i$ : *Pixels with depths in view i*

## Results



Grabcut [20]

Djelouah [8]

Ours

KinFu [19]

No Silh.

Ours