# Unsupervised Hierarchical Semantic Segmentation w/ Multiview Cosegmentation & Clustering Transformers
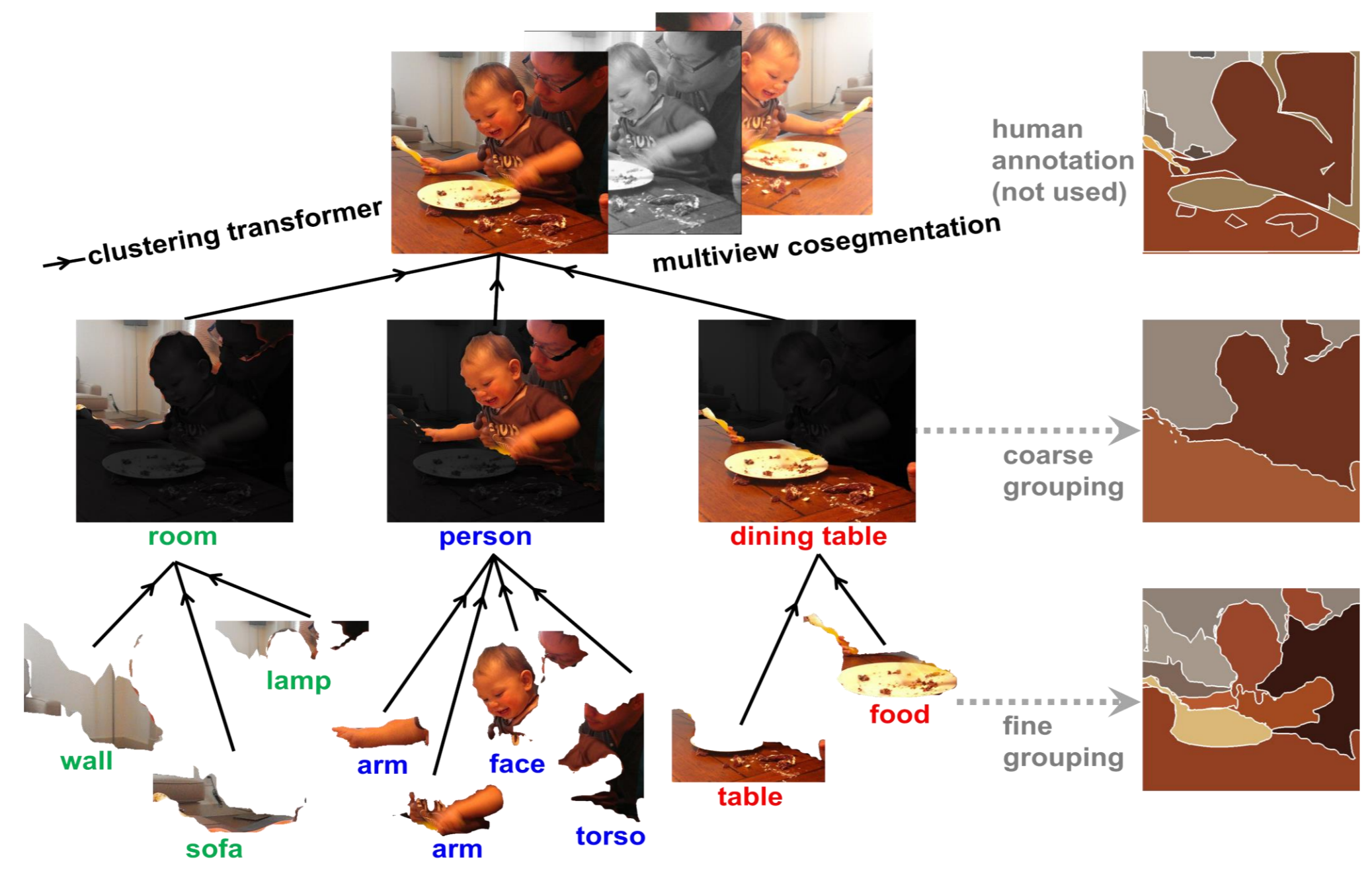
Tsung-Wei Ke    Jyh-Jing Hwang    Yunhui Guo    Xudong Wang    Stella X. Yu

## Unsupervised Hierarchical Segmentation



$$\text{NFCovering}(S' \to S_{fg}) = \frac{1}{|S_{fg}|} \sum_{R \in S_{fg}} \max_{R' \in S'} \frac{|R \cap R'|}{|R \cup R'|}$$



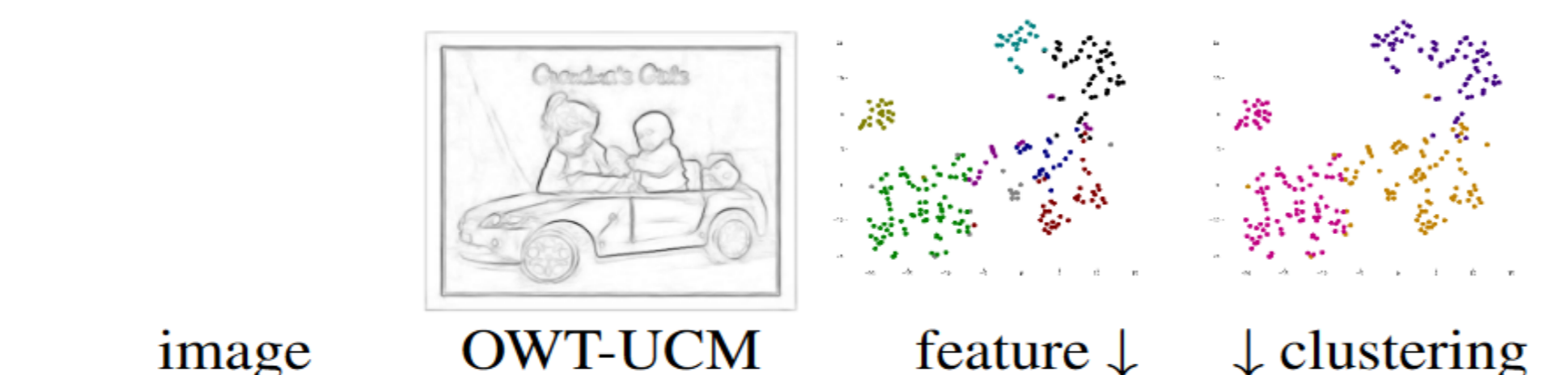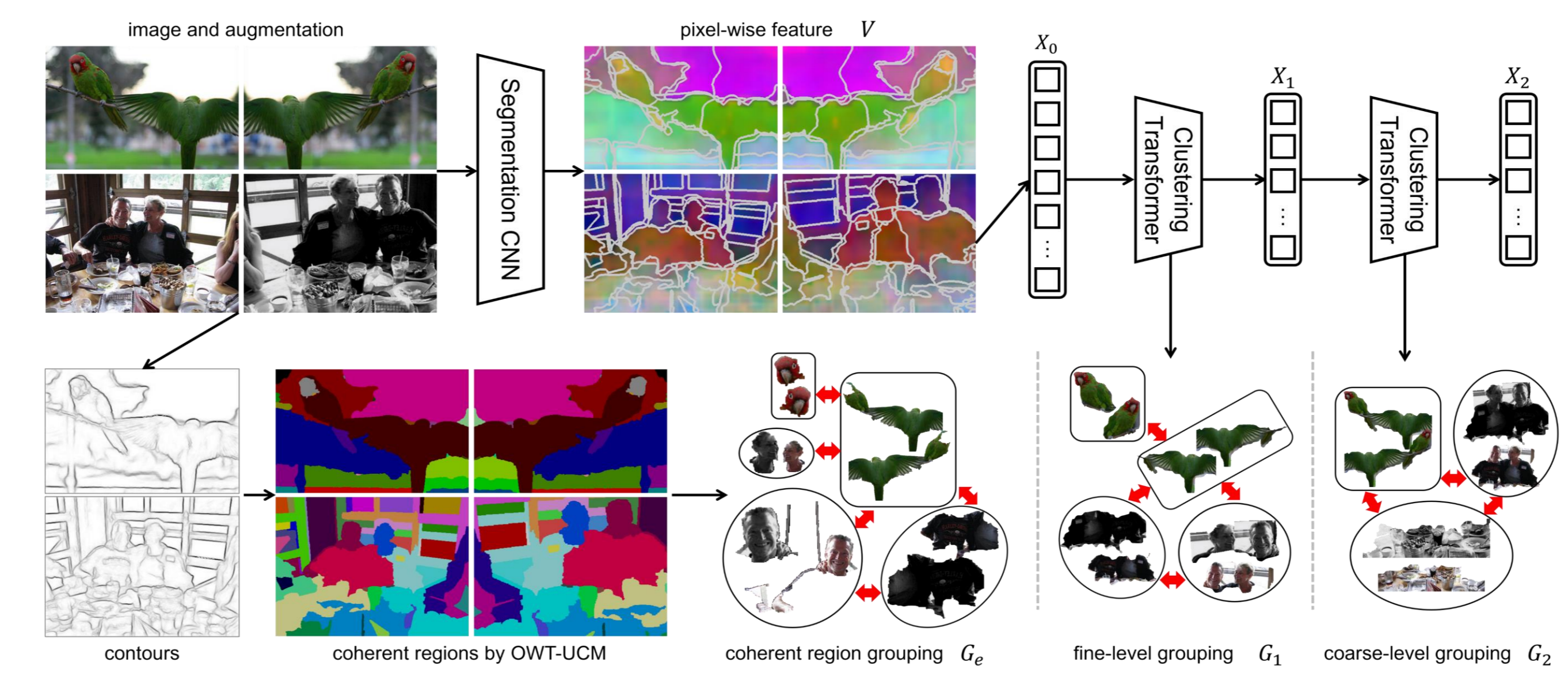VOC: Varying number of segmented regions

## Contributions



First unsupervised hierarchical semantic segmentation

First feature learning that embraces scale ambiguity

SOTA on unsupervised semantic segmentation

## Invariance: Multiview Cosegmentation

*Babies appear different but have the same semantics*



coherent region $G_e$    fine $G_1$    coarse $G_2$

## Hierarchical Segment Grouping Model



$$L(f) = \lambda_E L_f(G_e) + \lambda_F \sum_{l \geq 1} L_f(G_l) + \lambda_G L_g$$

Contrastive feature loss $L_f(G_e)$ grounds features by visual appearance

Contrastive feature loss $L_f(G_l)$ regularizes features by consistent hierarchy

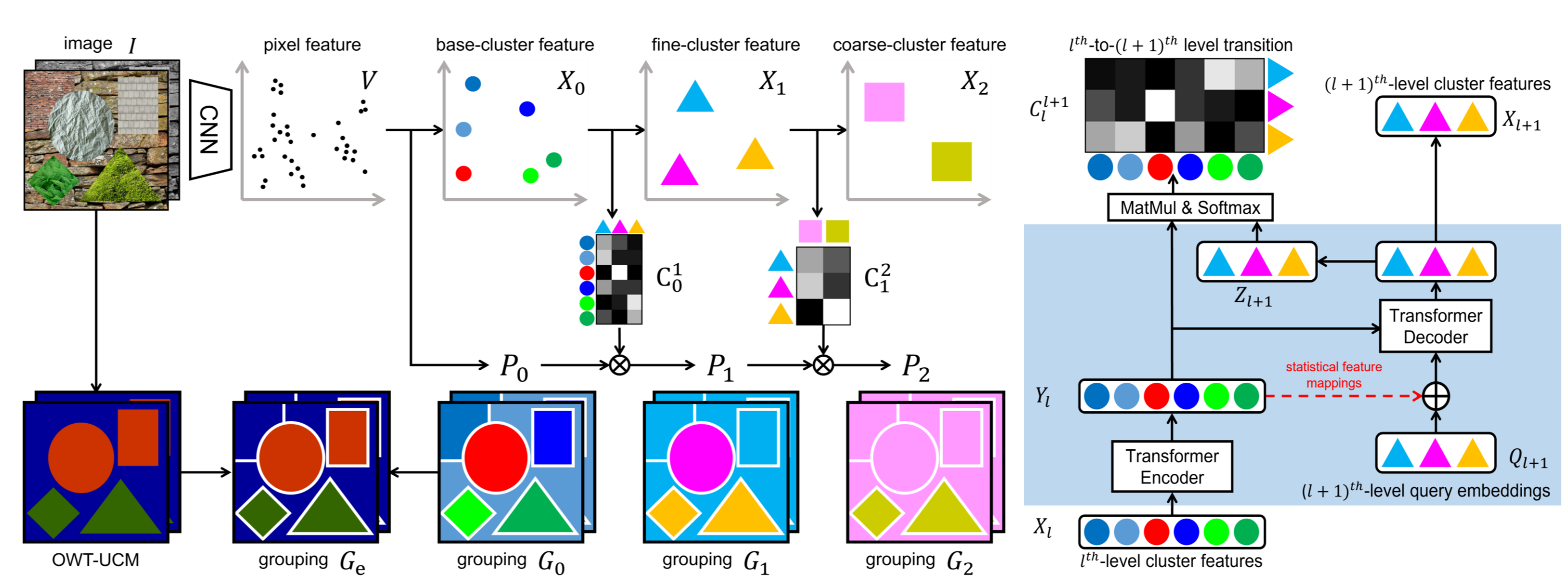Goodness of grouping $L_g$ desires balanced, compact, distinctive clusters

## Consistency: Clustering Transformers

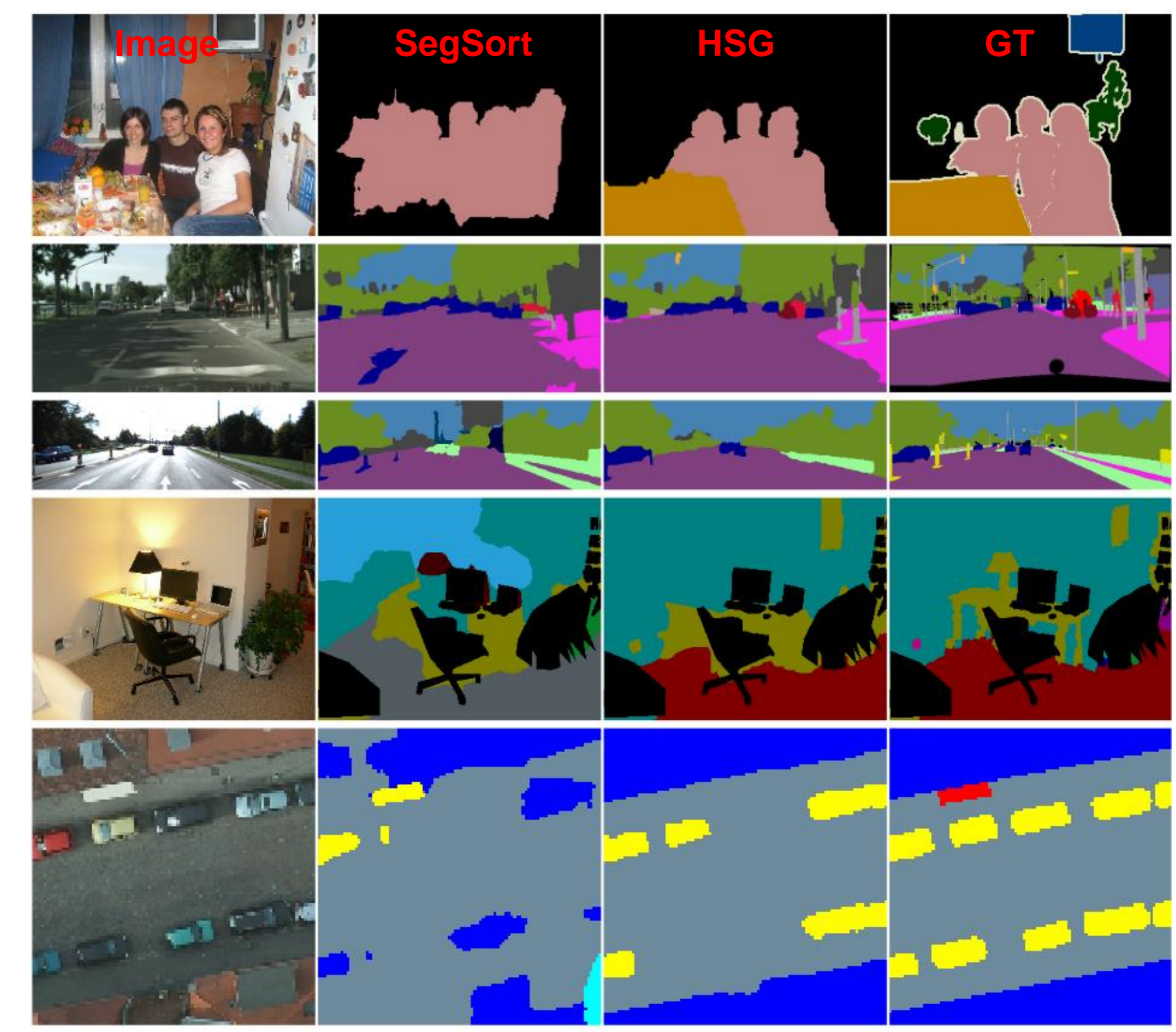*Face and body are parts of a whole in the visual scene*

Grouping Probability at Level $l$:    $P_l(a) = \text{Prob}(G_l = a|x)$

Transition Probability to Level $l+1$:    $C_l^{l+1}(a, b) = \text{Prob}(G_{l+1} = b | G_l = a)$

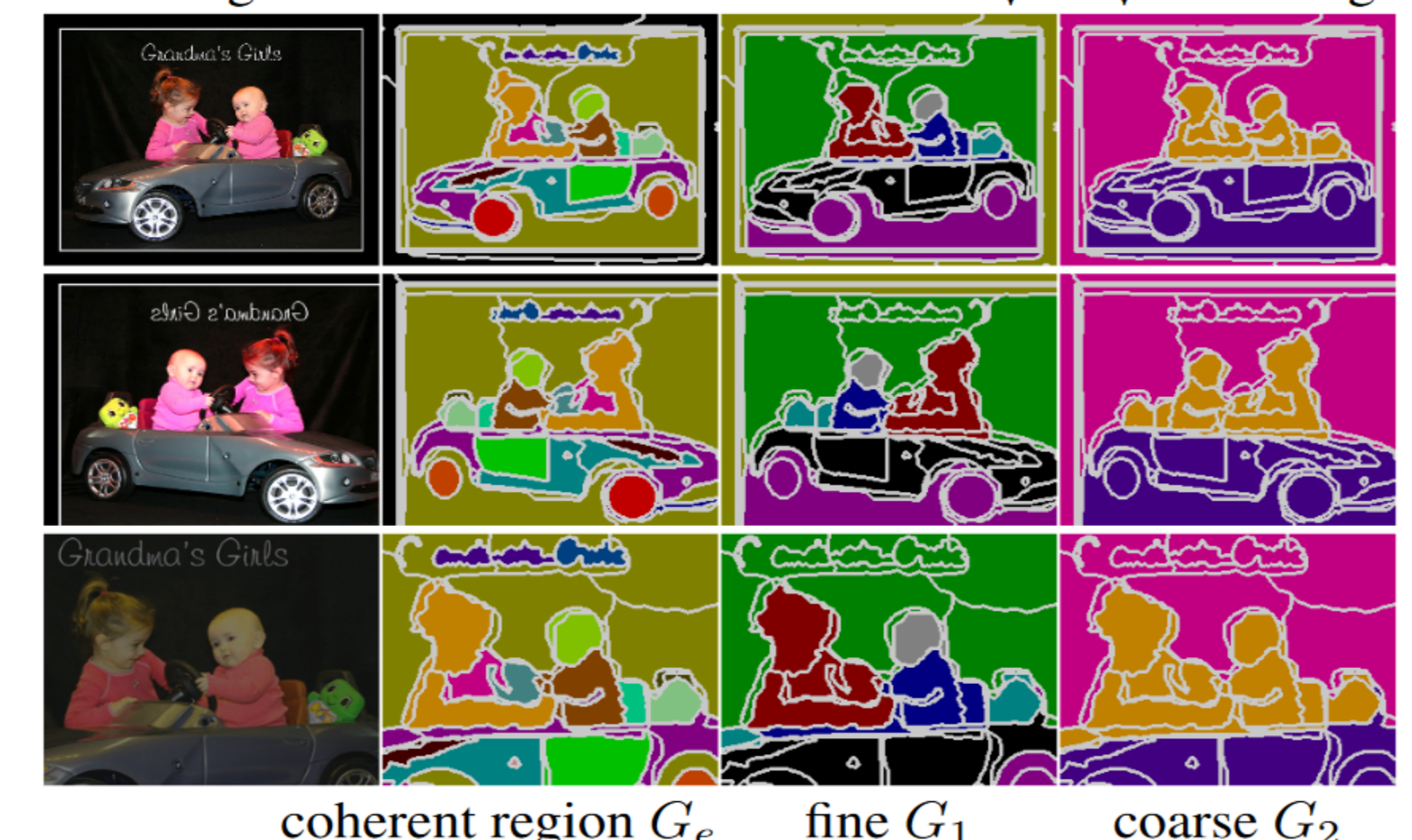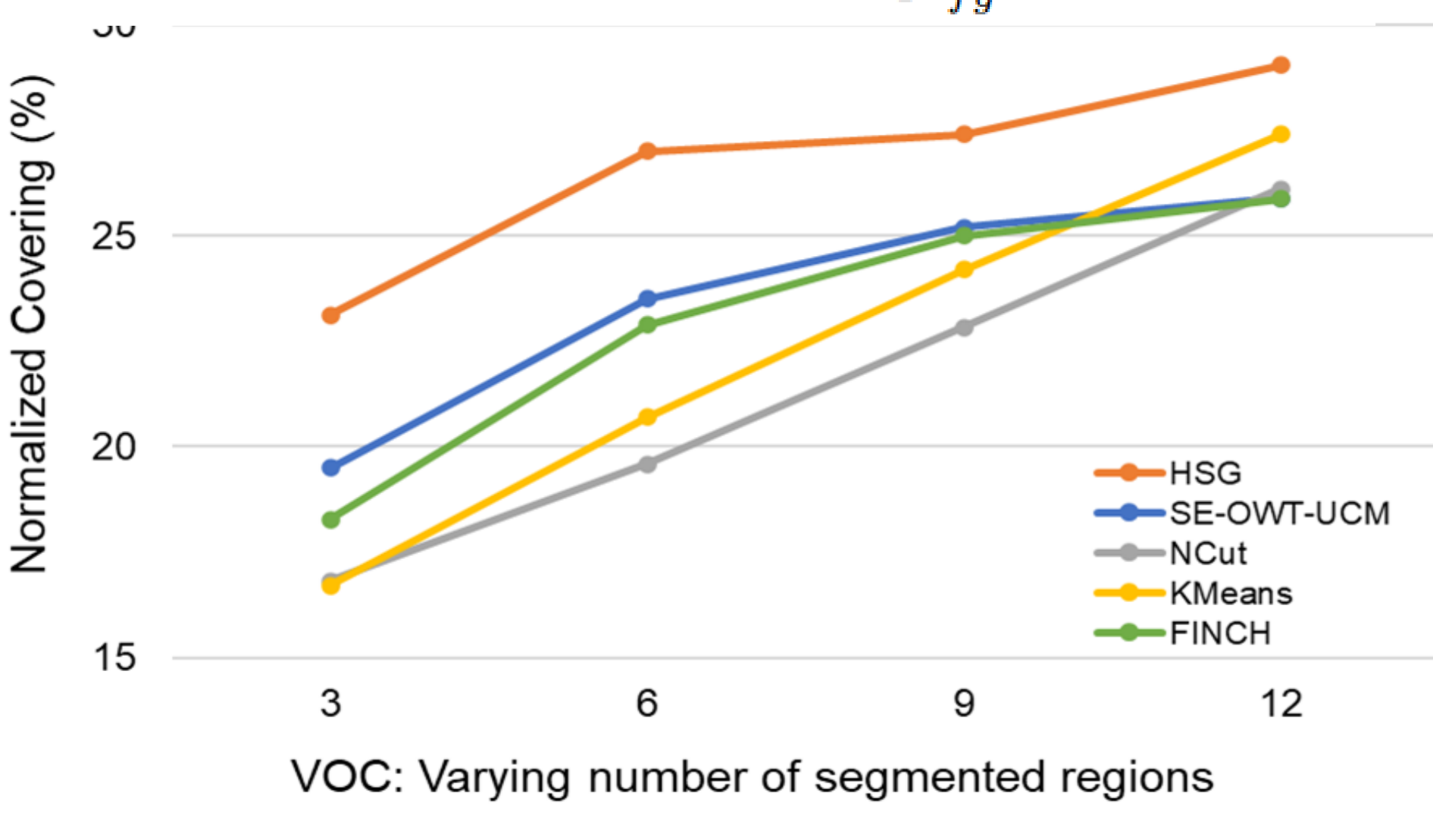Grouping Assignment at Level $l+1$: $P_{l+1} = P_l \times C_l^{l+1} = P_0 \times C_0^1 \times \cdots \times C_l^{l+1}$



## Unsupervised Semantic Segmentation



| Train | COCO | Cityscapes | KITTI |
|---|---|---|---|
| Test | VOC | Cityscapes | KITTI |
| Moco | 28.1 | 15.3 | 13.7 |
| DenseCL | 35.1 | 12.7 | 9.3 |
| Revisit | 35.1 | 17.1 | 17.0 |
| SegSort | 11.7 | 24.6 | 19.2 |
| Our HSG | **41.9** | **32.5** | **21.7** |

| Train & Test | COCO-stuff | Potsdam |
|---|---|---|
| DeepCluster | 19.9 | 29.2 |
| IIC | 27.7 | 45.4 |
| AC | 30.8 | 49.3 |
| SegSort | 49.9 | 59.0 |
| Our HSG | **57.6** | **67.4** |

## Unsupervised Contextual Retrievals