

## 1. Introduction:

➤ A core challenge in computer vision is to develop generative models of the world that capture rich contextual relationships among scene entities. Such models can serve different applications:

- **Scene Understanding:** regularizing the output of image descriptors in a Bayesian framework, generating sequences of unpredictable queries for testing computer vision systems (see the visual Turing test by Geman et al. [1]).
- **Robotics:** Simultaneous Localization and Mapping (SLAM), path planning, grasping and manipulating objects.
- **Computer Graphics:** creating synthetic content.

➤ Many man-made scenes are composed of multiple parallel supporting surfaces upon which instances from different object categories are placed [2].

➤ Designing and learning (from purely object-annotated images) 3D models which encode favored relationships but still accommodate real-world variability is not straightforward.

➤ We propose a new probabilistic, generative model of 3D scenes consisting of multiple objects lying on a plane. Our distribution is over random “**Generative Attributed Graphs (GAG)**” that encode favored layouts while accounting for variations in the number and relative poses of objects.

## 2. Proposed Model:

➤ A scene is described as a collection of object instances from different categories at different poses. Each object instance is associated with a vertex  $v \in V$  of a base graph  $g_0 \in G_0$  which captures contextual relationships among object instances.

➤ An attributed graph is a triple  $g = (g_0, c_V, \theta_V)$ , where  $c_V = \{c_v, v \in V\}$  and  $\theta_V = \{\theta_v, v \in V\}$  denote the set of category labels and 3D poses of objects, respectively.

➤ The model is a probability distribution on the space of attributed graphs conditioned on the environment's geometric properties  $T$ , specified by four sets of distributions:

1.  $p^{(0)}(n_{(0,1)}, \dots, n_{(0,K)}|T)$ : conditional joint distribution for the number of root nodes from each object category.
2.  $\{p^{(c)}(n_1, \dots, n_K), c \in C\}$ : the joint distribution of the number of children from each object category (Multi-type branching process). Restricted by a “Master Graph”.

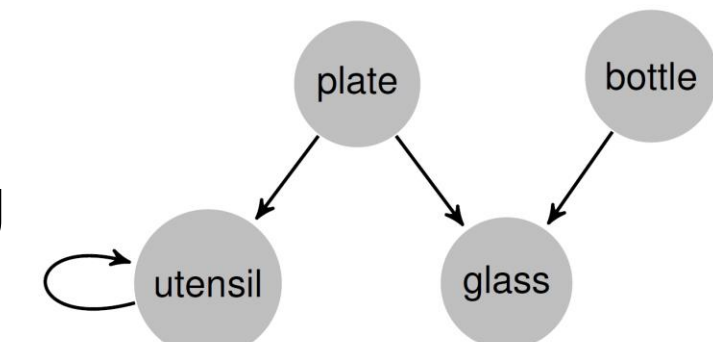
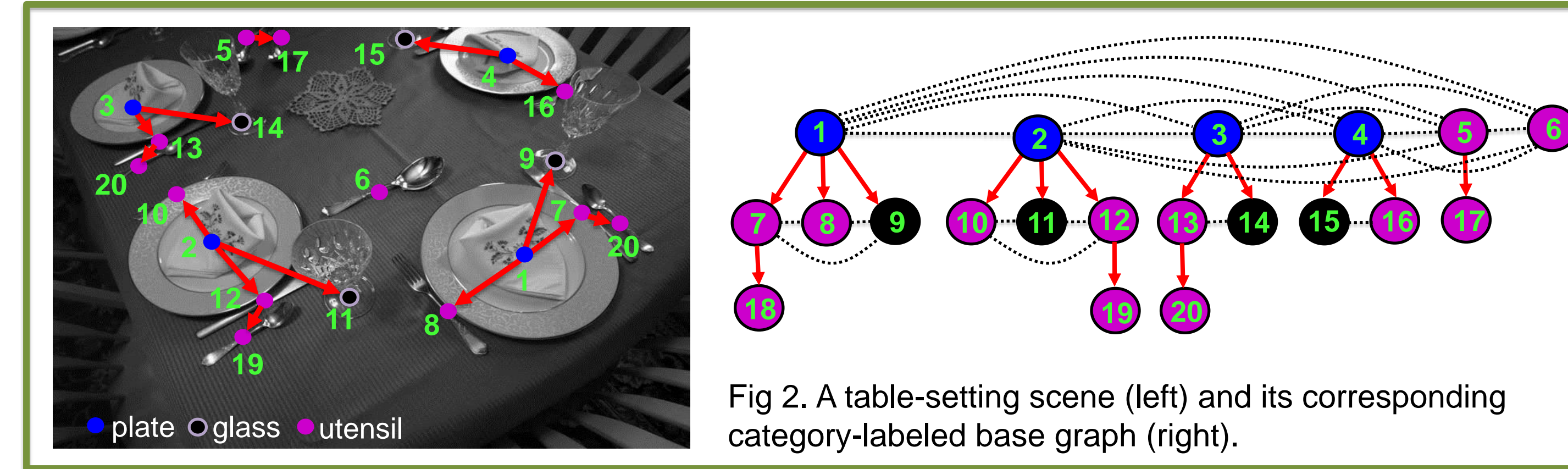


Fig 1. An Example “Master Graph”



3.  $p(\theta_{V_0}|c_{V_0}, T)$ : the joint distribution of the poses of the root nodes given  $T$ .
4.  $\{p(\theta_{ch(v)}|c_{ch(v)}, c_v, \theta_v, T), v \in V \setminus V_T\}$ : the joint distribution of the poses of the children of  $v$  given their parent's pose and the corresponding category labels and  $T$ .

The full distribution on attributed graphs  $g \in \mathcal{G}$ :

$$p(g|T) = p(g_0, c_V|T) \times p(\theta_V|g_0, c_V, T) \\ = p^{(0)}(n_{(0,1)}, \dots, n_{(0,K)}|T) p(\theta_{V_0}|c_{V_0}, T) \times \\ \prod_{v \in V \setminus V_T} p^{(c_v)}(n_{(v,1)}, \dots, n_{(v,K)}) p(\theta_{ch(v)}|c_{ch(v)}, c_v, \theta_v, T).$$

## 3. Model Learning:

➤ From annotated scenes:

- **Observable:**  $\mathcal{D} = \{c_V[j], \theta_V[j]\}_{j=1}^J$ , and **Hidden:**  $\mathcal{M} = \{g_0[j]\}_{j=1}^J$

➤ Parameter Estimation Using Expectation-Maximization (EM):

$$\Phi^{t+1} = \operatorname{argmax}_{\Phi} \sum_{j=1}^J \sum_{g_0[j]} p(g_0[j] | c_V[j], \theta_V[j], \Phi^t) \times \log p(g[j] | \Phi)$$

➤ Stochastic Expectation-Maximization using MCMC:

$$\Phi^{t+1} \approx \operatorname{argmax}_{\Phi} \sum_{j=1}^J \sum_{i=1}^N \log p(g^{(i)}[j] = (g_0^{(i)}[j], c_V[j], \theta_V[j]) | \Phi)$$

### Gibbs Sampling $p(g_0[j] | c_V[j], \theta_V[j], \Phi^t)$ :

➤ Conditional base graph distribution:

$$p(g_0[j] | c_V[j], \theta_V[j], \Phi^t) = \frac{1}{Z(c_V[j], \theta_V[j], \Phi^t)} \times p(g_0[j], c_V[j], \theta_V[j] | \Phi^t)$$

➤ For a scene with  $|V|$  annotated objects we can encode the base graph with:

$$\mathbf{z} = (z_1 = pa(v_1), z_2 = pa(v_2), \dots, z_{|V|} = pa(v_{|V|})), \quad g_0 = f(\mathbf{z})$$

**Step-1.** Begin with initial configuration  $i \leftarrow 1, l \leftarrow 1$ , and:

$$\mathbf{z}^{(0)} = (z_1^{(0)} = \emptyset, \dots, z_{|V|}^{(0)} = \emptyset)$$

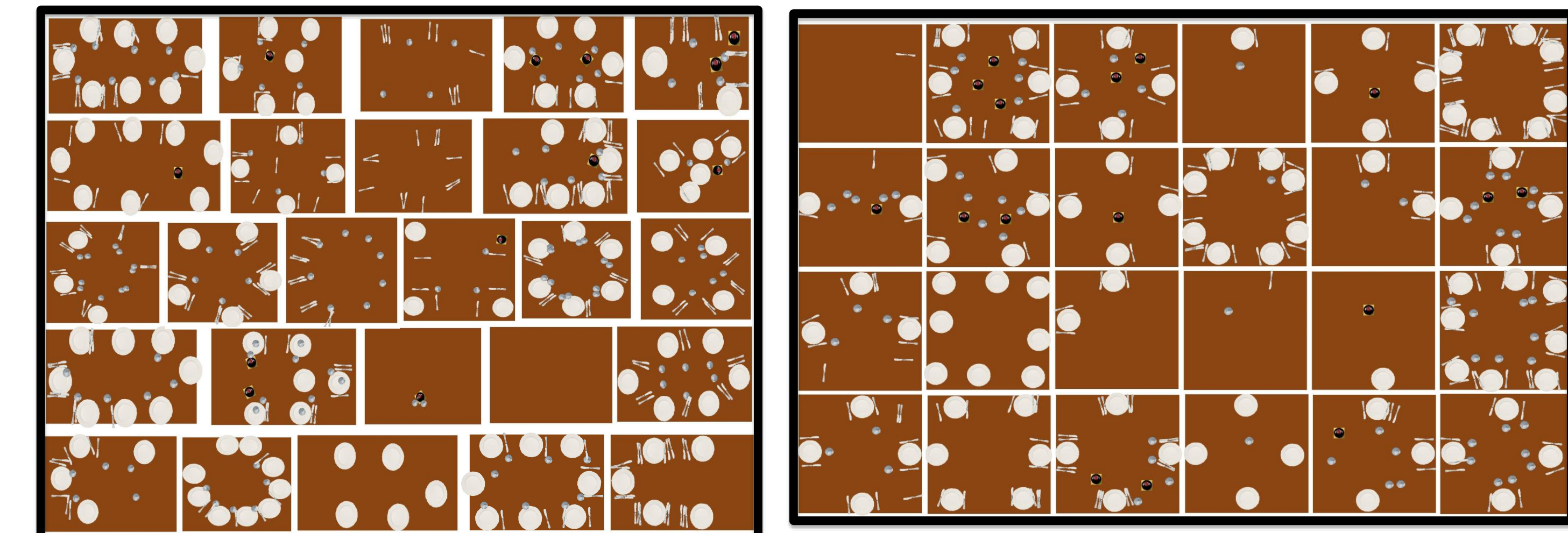
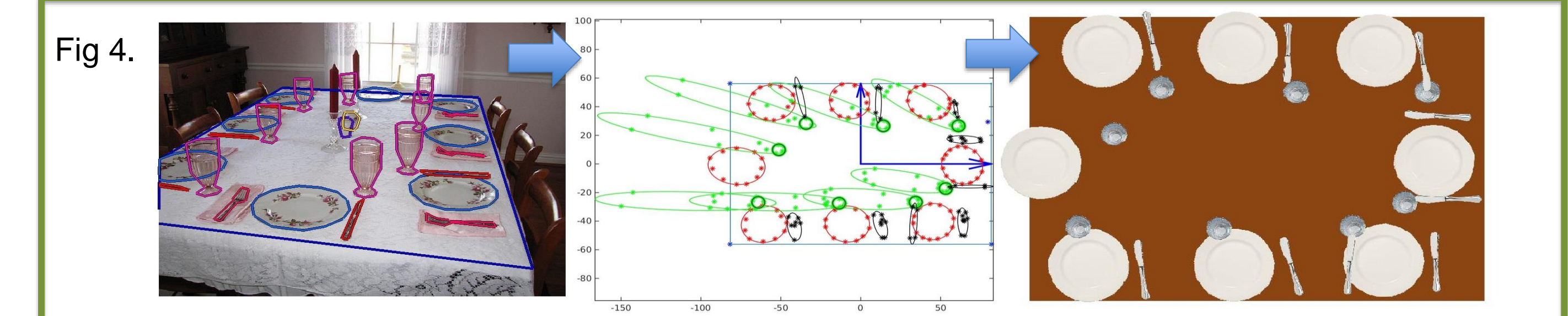
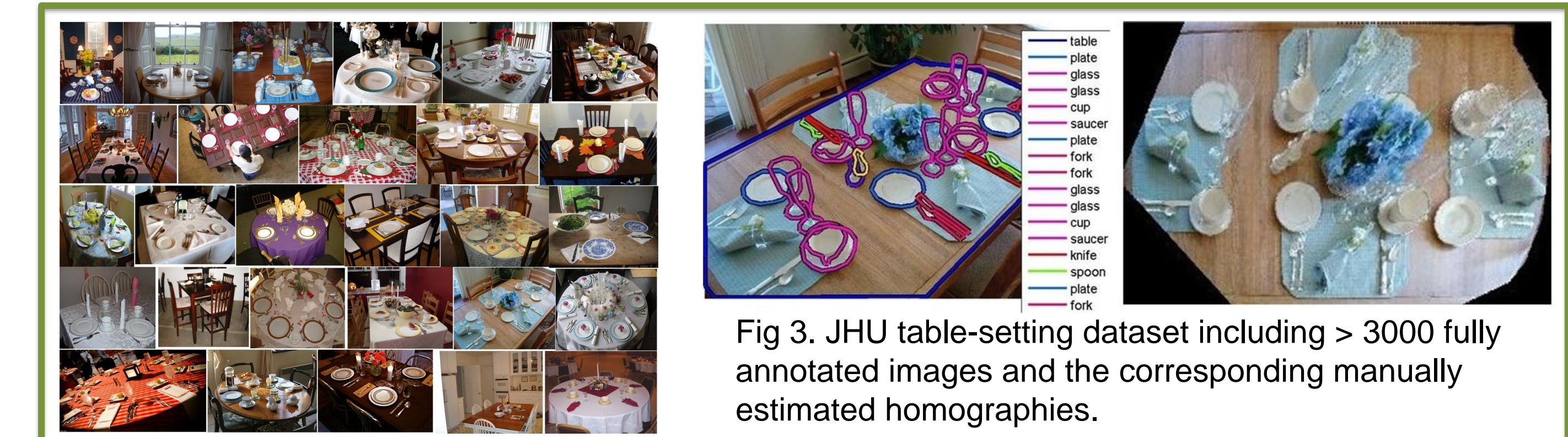
**Step-2.** Sweep  $\mathbf{z}$  by sampling one element at the time according to:

$$p(z_i) = \frac{p(g_0^{(l)}(z_i, c_V, \theta_V)}{\sum_{z_i \in S_{pa(v_i)}} p(g_0^{(l)}(z_i, c_V, \theta_V)}, \quad z_i \in S_{pa(v_i)}$$

**Step-3.** Generate the corresponding base graph sample  $g_0^{(l)}$ .

Set  $l \leftarrow l + 1$  and  $i \leftarrow (l \bmod |V|)$  and go back to Step-2.

## 4. Table-Setting Scenes:



**2D model from 3D model:**  $p(\xi_V, g, T, \mathcal{W}) = p(\xi_V|g, \mathcal{W}) p(g|T) p(\mathcal{W})$

### References:

- [1] D. Geman, et al. “A visual turing test for computer vision systems”. In PNAS, 2014.  
[2] S. Y. Bao, et al. “Toward coherent object detection and scene layout understanding”. In CVPR, 2010.