# EXPLAINING 3D OBJECT DETECTION THROUGH SHAPLEY VALUE-BASED ATTRIBUTION MAP

*Michihiro Kuroki and Toshihiko Yamasaki*

Dept. Information & Communication Engineering, The University of Tokyo, Tokyo

# Appendix

## Transformation of Equation

In this section, we describe the details of the transformation from Eq. 5 to Eq. 6 in the main paper. The expected value in Eq. 5 can be represented as follows:

$$\phi_i(f, \mathcal{X}) \approx \frac{1}{d} \sum_{k=1}^{d} \mathbb{E}\big[f(\mathcal{X}_{\mathbf{s}^k}) - f(\mathcal{X}_{\mathbf{s}^{k\prime}}) \mid \mathbf{s}_i^k = 1, \mathbf{s}_i^{k\prime} = 0\big], \tag{5}$$

$$= \frac{1}{d} \sum_{k=1}^{d} G_{\mathcal{X},k,i}. \tag{5a}$$

$$G_{\mathcal{X},k,i} = \mathbb{E}\big[f(\mathcal{X}_{\mathbf{s}^k}) - f(\mathcal{X}_{\mathbf{s}^{k\prime}}) \mid \mathbf{s}_i^k = 1, \mathbf{s}_i^{k\prime} = 0\big]. \tag{5b}$$

The expected value of Eq. 5b can be expressed as the summation of all combinations of two mask patterns. We denote two binary masks as $\mathbf{s}^{k,a}$ and $\mathbf{s}^{k,b}$, which exhibit patterns similar to that of $\mathbf{s}^k$. If duplication is permitted, we need to consider two conditions among the masks, namely $\mathbf{s}_i^{k,a} = 1, \mathbf{s}_i^{k,b} = 0$ and $\mathbf{s}_i^{k,a} = 0, \mathbf{s}_i^{k,b} = 1$.

$$\begin{aligned}
G_{\mathcal{X},k,i} = \sum_{\mathbf{m}^{k,a}} \sum_{\mathbf{m}^{k,b}} \Big\{ &\big(f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}})\big) P\big[\mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b} \mid \mathbf{s}_i^{k,a} = 1, \mathbf{s}_i^{k,b} = 0\big] \\
&+ \big(f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}})\big) P\big[\mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b} \mid \mathbf{s}_i^{k,a} = 0, , \mathbf{s}_i^{k,b} = 1\big] \Big\}.
\end{aligned} \tag{5c}$$

Here, $P$ denotes probability. This equation can be further transformed as follows:

$$\begin{aligned}
G_{\mathcal{X},k,i} = \sum_{\mathbf{m}^{k,a}} \sum_{\mathbf{m}^{k,b}} \Big\{ &\frac{(f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}})) P\big[\mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b}, \mathbf{s}_i^{k,a} = 1, \mathbf{s}_i^{k,b} = 0\big]}{P[\mathbf{s}_i^{k,a} = 1, \mathbf{s}_i^{k,b} = 0]} \\
&+ \frac{(f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}})) P\big[\mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b}, \mathbf{s}_i^{k,a} = 0, \mathbf{s}_i^{k,b} = 1\big]}{P[\mathbf{s}_i^{k,a} = 0, \mathbf{s}_i^{k,b} = 1]} \Big\}.
\end{aligned} \tag{5d}$$

$$\begin{aligned}
G_{\mathcal{X},k,i} = \frac{1}{P[\mathbf{s}_i^k = 1] \cdot P[\mathbf{s}_i^k = 0]} \sum_{\mathbf{m}^{k,a}} \sum_{\mathbf{m}^{k,b}} \Big\{ &\big(f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}})\big) P\big[\mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b}, \mathbf{s}_i^{k,a} = 1, \mathbf{s}_i^{k,b} = 0\big] \\
&+ \big(f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}})\big) P\big[\mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b}, \mathbf{s}_i^{k,a} = 1, \mathbf{s}_i^{k,b} = 0\big] \Big\},
\end{aligned} \tag{5e}$$

$$= \frac{1}{P[\mathbf{s}_i^k = 1] \cdot P[\mathbf{s}_i^k = 0]} \sum_{\mathbf{m}^{k,a}} \sum_{\mathbf{m}^{k,b}} \left\{ \left( f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}}) \right) \left( \mathbf{m}_i^{k,a} - \mathbf{m}_i^{k,b} \right) P\left[ \mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b} \right] \right\}. \tag{5f}$$

We now aim to reformulate the summation over $\mathbf{m}^{k,b}$ in terms of its expected value.

$$\sum_{\mathbf{m}^{k,a}} \sum_{\mathbf{m}^{k,b}} \left\{ \left( f(\mathcal{X}_{\mathbf{m}^{k,a}}) - f(\mathcal{X}_{\mathbf{m}^{k,b}}) \right) \left( \mathbf{m}_i^{k,a} - \mathbf{m}_i^{k,b} \right) P\left[ \mathbf{s}^{k,a} = \mathbf{m}^{k,a}, \mathbf{s}^{k,b} = \mathbf{m}^{k,b} \right] \right\},$$

$$= \sum_{\mathbf{m}^{k,a}} \left\{ f(\mathcal{X}_{\mathbf{m}^{k,a}}) \cdot \mathbf{m}_i^{k,a} - f(\mathcal{X}_{\mathbf{m}^{k,a}}) \cdot \mathbb{E}\left[ \mathbf{s}_i^{k,b} \right] - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^{k,b}}) \right] \cdot \mathbf{m}_i^{k,a} + \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^{k,b}}) \cdot \mathbf{s}_i^{k,b} \right] \right\} P\left[ \mathbf{s}^{k,a} = \mathbf{m}^{k,a} \right], \tag{5g}$$

$$\approx \sum_{\mathbf{m}^{k,a}} \left\{ f(\mathcal{X}_{\mathbf{m}^{k,a}}) \cdot \mathbf{m}_i^{k,a} - f(\mathcal{X}_{\mathbf{m}^{k,a}}) \cdot \mathbb{E}\left[ \mathbf{s}_i^{k,b} \right] - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^{k,b}}) \right] \cdot \mathbf{m}_i^{k,a} + \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^{k,b}}) \right] \cdot \mathbb{E}\left[ \mathbf{s}_i^{k,b} \right] \right\} P\left[ \mathbf{s}^{k,a} = \mathbf{m}^{k,a} \right], \tag{5h}$$

$$= \sum_{\mathbf{m}^{k,a}} \left\{ \left( f(\mathcal{X}_{\mathbf{m}^{k,a}}) - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^{k,b}}) \right] \right) \left( \mathbf{m}_i^{k,a} - \mathbb{E}\left[ \mathbf{s}_i^{k,b} \right] \right) \right\} P\left[ \mathbf{s}^{k,a} = \mathbf{m}^{k,a} \right]. \tag{5i}$$

In the transformation, we assumed independence between $f(\mathcal{X}_{\mathbf{s}^{k,b}})$ and $\mathbf{s}_i^{k,b}$. Given that $\mathbf{m}^{k,a}$ and $\mathbf{m}^{k,b}$ follow the same distribution of $\mathbf{s}^k$, we can rewrite Eq. 5i as follows.

$$G_{\mathcal{X},k,i} \approx \frac{1}{P[\mathbf{s}_i^k = 1] \cdot P[\mathbf{s}_i^k = 0]} \sum_{\mathbf{m}^k} \left\{ \left( f(\mathcal{X}_{\mathbf{m}^k}) - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^k}) \right] \right) \left( \mathbf{m}_i^k - \mathbb{E}\left[ \mathbf{s}_i^k \right] \right) \right\} P\left[ \mathbf{s}^k = \mathbf{m}^k \right], \tag{5j}$$

$$= \frac{1}{\mathbb{E}\left[ \mathbf{s}_i^k \right] \left( 1 - \mathbb{E}\left[ \mathbf{s}_i^k \right] \right)} \sum_{\mathbf{m}^k} \left\{ \left( f(\mathcal{X}_{\mathbf{m}^k}) - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^k}) \right] \right) \left( \mathbf{m}_i^k - \mathbb{E}\left[ \mathbf{s}_i^k \right] \right) \right\} P\left[ \mathbf{s}^k = \mathbf{m}^k \right]. \tag{5k}$$

Using the definition of covariance, we ultimately rewrite the summation as the expected values over $\mathbf{s}^k$.

$$G_{\mathcal{X},k,i} \approx \frac{\mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^k}) \cdot \mathbf{s}_i^k \right] - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^k}) \right] \cdot \mathbb{E}\left[ \mathbf{s}_i^k \right]}{\mathbb{E}\left[ \mathbf{s}_i^k \right] \left( 1 - \mathbb{E}\left[ \mathbf{s}_i^k \right] \right)}. \tag{5l}$$

$$\phi_i(f, \mathcal{X}) \approx \frac{1}{d} \sum_{k=1}^{d} \frac{\mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^k}) \cdot \mathbf{s}_i^k \right] - \mathbb{E}\left[ f(\mathcal{X}_{\mathbf{s}^k}) \right] \cdot \mathbb{E}\left[ \mathbf{s}_i^k \right]}{\mathbb{E}\left[ \mathbf{s}_i^k \right] \cdot \left( 1 - \mathbb{E}\left[ \mathbf{s}_i^k \right] \right)} \tag{6}$$

## Pseudocode

The pseudocode describing our method is shown in Algorithm 1.

---

**Algorithm 1** Pseudocode for computing attribution map $\Phi$

---

**Inputs:** The number of samplings $N$, number of approximation layers $L$, object detector function $F$, input point cloud $\mathcal{X}$, explanation target detection $\mathcal{D}_t$, detection score function $Sim(\cdot)$, and all-ones mask $\mathbf{1}$.

**Outputs:** Attribution map $\Phi$

1: $\Phi \leftarrow O$
2: **for** $l = 1, \ldots, L$ **do**
3:      $\Phi^l \leftarrow O$
4:      sum_score $\leftarrow 0$, sum_mask $\leftarrow O$, sum_score_mask $\leftarrow O$
5:      **for** $r = 1, \ldots, N$ **do**
6:          $\mathbf{s}^{l_r} \leftarrow$ The input point cloud space is divided into voxel units. The voxels are selected randomly with probability $p = \frac{l}{L+1}$, and a point $i$ within the unselected voxels is masked (i.e. $\mathbf{s}_i^{l_r} = 0$).
7:          $f\left(\mathcal{X}_{\mathbf{s}^{l_r}}\right) \leftarrow \max_{\mathcal{D}_j \in F(\mathcal{X}_{\mathbf{s}^{l_r}})} Sim(\mathcal{D}_t, \mathcal{D}_j)$
8:          sum_score $\leftarrow$ sum_score $+ f\left(\mathcal{X}_{\mathbf{s}^{l_r}}\right)$
9:          sum_mask $\leftarrow$ sum_mask $+ \mathbf{s}^{l_r}$
10:         sum_score_mask $\leftarrow$ sum_score_mask $+ f\left(\mathcal{X}_{\mathbf{s}^{l_r}}\right) \cdot \mathbf{s}^{l_r}$
11:      **end for**
12:      $\overline{f\left(\mathcal{X}_{\mathbf{s}^l}\right)} \leftarrow$ sum_score$/N$
13:      $\overline{\mathbf{s}^l} \leftarrow$ sum_mask$/N$
14:      $\overline{f\left(\mathcal{X}_{\mathbf{s}^l}\right) \cdot \mathbf{s}^l} \leftarrow$ sum_score_mask$/N$
15:      $\Phi_l \leftarrow \frac{1}{L} \cdot \left\{ \overline{f\left(\mathcal{X}_{\mathbf{s}^l}\right) \cdot \mathbf{s}^l} - \overline{f\left(\mathcal{X}_{\mathbf{s}^l}\right)} \cdot \overline{\mathbf{s}^l} \right\} \oslash \left\{ \overline{\mathbf{s}^l} \odot \left(\mathbf{1} - \overline{\mathbf{s}^l}\right) \right\}$
16:      $\Phi \leftarrow \Phi + \Phi_l$
17: **end for**
18: **return** $\Phi$

---