# $E(2)$-Equivariant Vision Transformer
# (Supplementary Material)

**Renjun Xu**[*1]          **Kaifan Yang**[*1,2]          **Ke Liu**[*†1,2]          **Fengxiang He**[3,2]

[1]College of Computer Science and Technology, Zhejiang University
[2]JD Explore Academy, JD.com, Inc.
[3]AIAI, School of Informatics, University of Edinburgh

## A   ERRORS IN GSA-NETS

In this part, we first review the proof of equivariance of the GSA-Nets [Romero and Cordonnier, 2020], and then point out the mistakes in the proof process using the positional encoding as:

$$\rho((i,\tilde{h}),(j,\hat{h})) = \rho^P(x(j) - x(i), \tilde{h}^{-1}\hat{h})$$

### A.1   DEFINITIONS AND NOTATIONS.

#### A.1.1   Definition of Group Equivariant Self-Attention.

If the group self-attention formulation $m_{\mathcal{G}}^r[f,\rho](i,\hbar)$ is $\mathcal{G}$-equivariant, if and only if it satisfies:

$$m_{\mathcal{G}}^r[\mathcal{L}_g[f],\rho](i,\hbar) = \mathcal{L}_g[m_{\mathcal{G}}^r[f,\rho]](i,\hbar), \quad g \in \mathcal{G}$$

#### A.1.2   Input under $g$-Transformed

A $g$-transformed input can be expressed as:

$$\mathcal{L}_g[f](i,\tilde{h}) = \mathcal{L}_y\mathcal{L}_{\bar{h}}[f](i,\tilde{h}) = f(\rho^{-1}(\bar{h}^{-1}(\rho(i) - y)), \bar{h}^{-1}\tilde{h}),$$
$$g = (y,\bar{h}),\ y \in \mathbb{R}^d,\ \bar{h} \in \mathcal{H}.$$

## B   PROOF OF GE-VIT

In this section, we prove that GE-ViT is group equivariant. For brevity, we also use the substitutions:

$$\bar{i} = x^{-1}(\bar{h}^{-1}(x(i) - y)) \Rightarrow i = x^{-1}(\bar{h}x(\bar{i}) + y)), \tilde{h}' = \bar{h}^{-1}\tilde{h}$$

and

$$\bar{j} = x^{-1}(\bar{h}^{-1}(x(j) - y)) \Rightarrow j = x^{-1}(\bar{h}x(\bar{j}) + y)), \hat{h}' = \bar{h}^{-1}\hat{h}$$

---

[*]Contributed equally.
[†]Corresponding author: Ke Liu

The complete proof process is as follows:

$$m_{\mathcal{G}}^r\big[\mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[f],\rho\big](i,\hbar) \tag{1}$$

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\tilde{\hbar}\in\mathcal{H}}\sum_{(j,\hat{\hbar})\in n(i,\tilde{\hbar})}\sigma_{j,\hat{\hbar}}\big(\langle\varphi_{\text{qry}}^{(h)}(\mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[f](i,\tilde{\hbar})),\varphi_{\text{key}}^{(h)}(\mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[f](j,\hat{\hbar})$$

$$+ \mathcal{L}_{\hbar}[\rho]((i,\tilde{\hbar}),(j,\hat{\hbar}))\rangle\varphi_{\text{val}}^{(h)}(\mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[f](j,\hat{\hbar}))\Big)$$

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\tilde{\hbar}\in\mathcal{H}}\sum_{(j,\hat{\hbar})\in n(i,\tilde{\hbar})}\sigma_{j,\hat{\hbar}}\big(\langle\varphi_{\text{qry}}^{(h)}(f(x^{-1}(\bar{\hbar}^{-1}(x(i)-y)),\bar{\hbar}^{-1}\tilde{\hbar})), \tag{2}$$

$$\varphi_{\text{key}}^{(h)}(f(x^{-1}(\bar{\hbar}^{-1}(x(j)-y)),\bar{\hbar}^{-1}\hat{\hbar})+\mathcal{L}_{\hbar}[\rho]((i,\tilde{\hbar}),(j,\hat{\hbar}))\rangle$$

$$\varphi_{\text{val}}^{(h)}(f(x^{-1}(\bar{\hbar}^{-1}(x(j)-y)),\bar{\hbar}^{-1}\hat{\hbar}))\Big)$$

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\bar{\hbar}\tilde{\hbar}'\in\mathcal{H}}\sum_{(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}')}\sigma_{x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}'}\big(\langle\varphi_{\text{qry}}^{(h)}(f(\bar{i},\tilde{\hbar}')),\varphi_{\text{key}}^{(h)}(f(\bar{j},\hat{\hbar}') \tag{3}$$

$$+ \mathcal{L}_{\hbar}[\rho]((x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}'),(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}'))\rangle\varphi_{\text{val}}^{(h)}(f(\bar{j},\hat{\hbar}'))\Big)$$

By using the definition:

$$\rho((i,\tilde{\hbar}),(j,\hat{\hbar})) = \rho^P(x(j)-x(i),\tilde{\hbar}\hat{\hbar}^{-1}\tilde{\hbar})$$

and

$$\mathcal{L}_{\hbar}[\rho]((i,\tilde{\hbar}),(j,\hat{\hbar})) = \rho^P(\hbar^{-1}(x(j)-x(i)),\hbar^{-1}(\tilde{\hbar}\hat{\hbar}^{-1}\tilde{\hbar})).$$

The above formula can be further derived:

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\bar{\hbar}\tilde{\hbar}'\in\mathcal{H}}\sum_{(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}')}\sigma_{x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}'}\big(\langle\varphi_{\text{qry}}^{(h)}(f(\bar{i},\tilde{\hbar}')),\varphi_{\text{key}}^{(h)}(f(\bar{j},\hat{\hbar}') \tag{4}$$

$$+ \rho^P(\hbar^{-1}(\bar{\hbar}x(\bar{j})+y-(\bar{\hbar}x(\bar{i})+y)),\hbar^{-1}(\bar{\hbar}\tilde{\hbar}')(\bar{\hbar}\hat{\hbar}')^{-1}(\bar{\hbar}\tilde{\hbar}'))\rangle\varphi_{\text{val}}^{(h)}(f(\bar{j},\hat{\hbar}'))\Big)$$

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\bar{\hbar}\tilde{\hbar}'\in\mathcal{H}}\sum_{(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}')}\sigma_{x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}'}\big(\langle\varphi_{\text{qry}}^{(h)}(f(\bar{i},\tilde{\hbar}')),\varphi_{\text{key}}^{(h)}(f(\bar{j},\hat{\hbar}') \tag{5}$$

$$+ \rho^P(\hbar^{-1}(\bar{\hbar}x(\bar{j})+y-(\bar{\hbar}x(\bar{i})+y)),\hbar^{-1}\bar{\hbar}\tilde{\hbar}'\hat{\hbar}'^{-1}\tilde{\hbar}')\rangle\varphi_{\text{val}}^{(h)}(f(\bar{j},\hat{\hbar}'))\Big)$$

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\bar{\hbar}\tilde{\hbar}'\in\mathcal{H}}\sum_{(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}')}\sigma_{x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}'}\big(\langle\varphi_{\text{qry}}^{(h)}(f(\bar{i},\tilde{\hbar}')),\varphi_{\text{key}}^{(h)}(f(\bar{j},\hat{\hbar}') \tag{6}$$

$$+ \rho^P(\hbar^{-1}\bar{\hbar}(x(\bar{j})-x(\bar{i})),\tilde{\hbar}'\hat{\hbar}'^{-1}\tilde{\hbar}'))\rangle\varphi_{\text{val}}^{(h)}(f(\bar{j},\hat{\hbar}'))\Big)$$

$$= \varphi_{\text{out}}\Big(\bigcup_{h\in[H]}\sum_{\bar{\hbar}\tilde{\hbar}'\in\mathcal{H}}\sum_{(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}')}\sigma_{x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}'}\big(\langle\varphi_{\text{qry}}^{(h)}(f(\bar{i},\tilde{\hbar}')),\varphi_{\text{key}}^{(h)}(f(\bar{j},\hat{\hbar}') \tag{7}$$

$$+ \mathcal{L}_{\bar{\hbar}^{-1}\hbar}[\rho]((\bar{i},\tilde{\hbar}'),(\bar{j},\hat{\hbar}'))\rangle\varphi_{\text{val}}^{(h)}(f(\bar{j},\hat{\hbar}'))\Big)$$

The subsequent proof is similar to the GSA-Nets [Romero and Cordonnier, 2020]. For unimodular groups, the area of summation remains equal for any transformation $g\in\mathcal{G}$, which means that:

$$\sum_{(x^{-1}(\bar{\hbar}x(\bar{j})+y),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})+y),\bar{\hbar}\tilde{\hbar}')}[\cdot] \quad = \quad \sum_{(x^{-1}(\bar{\hbar}x(\bar{j})),\bar{\hbar}\hat{\hbar}')\in n(x^{-1}(\bar{\hbar}x(\bar{i})),\bar{\hbar}\tilde{\hbar}')}[\cdot]$$

$$= \quad \sum_{(x^{-1}(x(\bar{j})),\hat{\hbar}')\in n(x^{-1}(x(\bar{i})),\tilde{\hbar}')}[\cdot]$$

$$= \quad \sum_{(\bar{j},\hat{\hbar}')\in n(\bar{i},\tilde{\hbar}')}[\cdot].$$

and because of the basic properties of groups, we can get $\sum_{\bar{\hbar}\tilde{\hbar}'\in\mathcal{H}}[\cdot]=\sum_{\tilde{\hbar}'\in\mathcal{H}}[\cdot]$. Consequently, the above formula can be further simplified as:

$$
\begin{aligned}
m_{\mathcal{G}}^r\big[\mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[f],\rho\big](i,\hbar) =& \varphi_{\text{out}}\Big( \bigcup_{h\in[H]} \sum_{\tilde{\hbar}'\in\mathcal{H}} \sum_{(\bar{j},\hat{\hbar}')\in\bar{n}(\bar{i},\tilde{\hbar}')} \sigma_{\bar{j},\hat{\hbar}'}\big(\langle\varphi_{\text{qry}}^{(h)}(f(\bar{i},\tilde{\hbar}')), \\
& \varphi_{\text{key}}^{(h)}(f(\bar{j},\hat{\hbar}') + \mathcal{L}_{\bar{\hbar}^{-1}\hbar}[\rho]((\bar{i},\tilde{\hbar}'),(\bar{j},\hat{\hbar}')))\rangle\big)\varphi_{\text{val}}^{(h)}(f(\bar{j},\hat{\hbar}'))\Big) \\
=& m_{\mathcal{G}}^r[f,\rho](\bar{i},\bar{\hbar}^{-1}\hbar) \\
=& m_{\mathcal{G}}^r[f,\rho](x^{-1}(\bar{\hbar}^{-1}(x(i)-y)),\bar{\hbar}^{-1}\hbar) \\
=& \mathcal{L}_y\mathcal{L}_{\bar{\hbar}}\big[m_{\mathcal{G}}^r[f,\rho]\big](i,\hbar).
\end{aligned}
\tag{8}
$$

From the above formula, it can be seen that:

$$
m_{\mathcal{G}}^r[\mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[f],\rho](i,\hbar) = \mathcal{L}_y\mathcal{L}_{\bar{\hbar}}[m_{\mathcal{G}}^r[f,\rho]](i,\hbar),
$$

which is the same as:

$$
m_{\mathcal{G}}^r[\mathcal{L}_g[f],\rho](i,\hbar) = \mathcal{L}_g[m_{\mathcal{G}}^r[f,\rho]](i,\hbar), \quad g\in\mathcal{G}.
$$

Therefore, with the positional encoding we proposed:

$$
\rho((i,\tilde{\hbar}),(j,\hat{\hbar})) = \rho^P(x(j)-x(i),\tilde{\hbar}\hat{\hbar}^{-1}\tilde{\hbar}),
$$

the group self-attention is group equivariant.

## References

David W Romero and Jean-Baptiste Cordonnier. Group equivariant stand-alone self-attention for vision. In *International Conference on Learning Representations*, 2020.