



Semi-supervised Hand Appearance Recovery via Structure Disentanglement and Dual Adversarial Discrimination

Zimeng Zhao, Binghui Zuo, Zhiyu Long and Yangang Wang

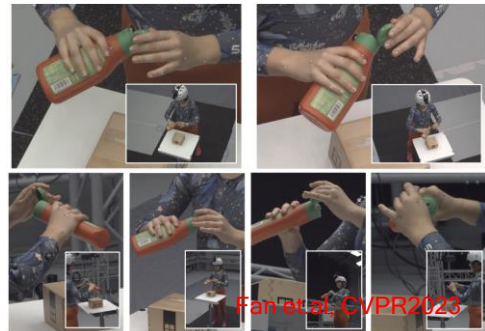
Southeast University

WED-AM-372

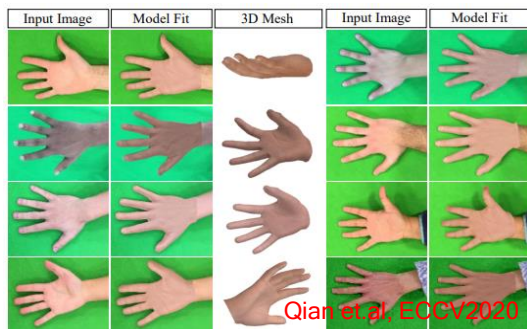
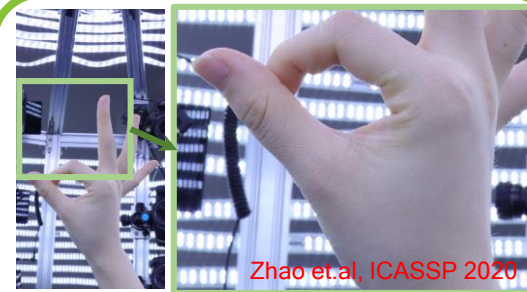


Human Hand Data Capture

Separate Acquisition

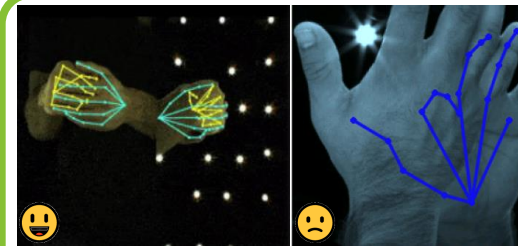


Accurate Pose:
Marker-based Capture

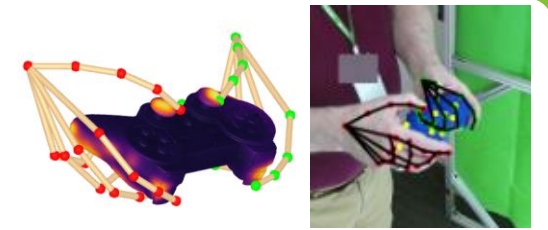


Realistic Appearance:
Bare Skin Capture

Synchronous Acquisition



Moon et.al, ECCV 2020



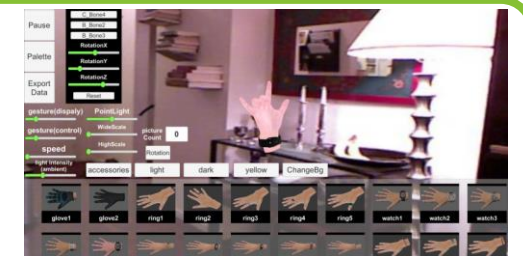
Brahmbhatt et.al, ECCV 2020

Marker-less Procedure with Pose Estimator



Hasson et.al, CVPR 2022

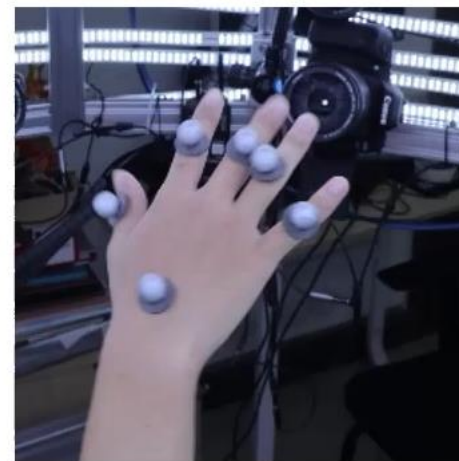
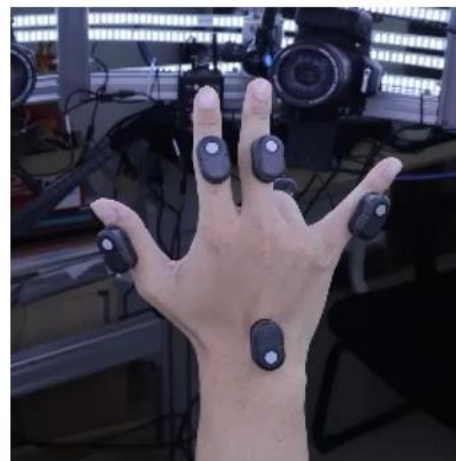
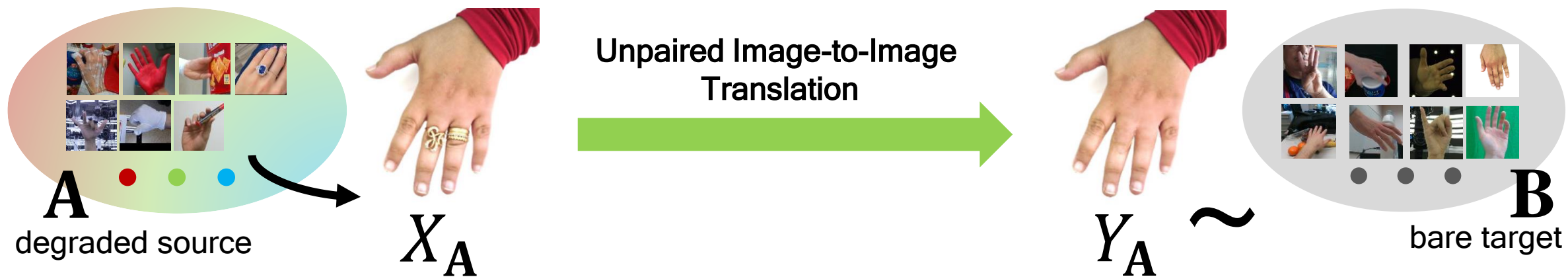
Synthetic Procedure with Differentiable Rendering



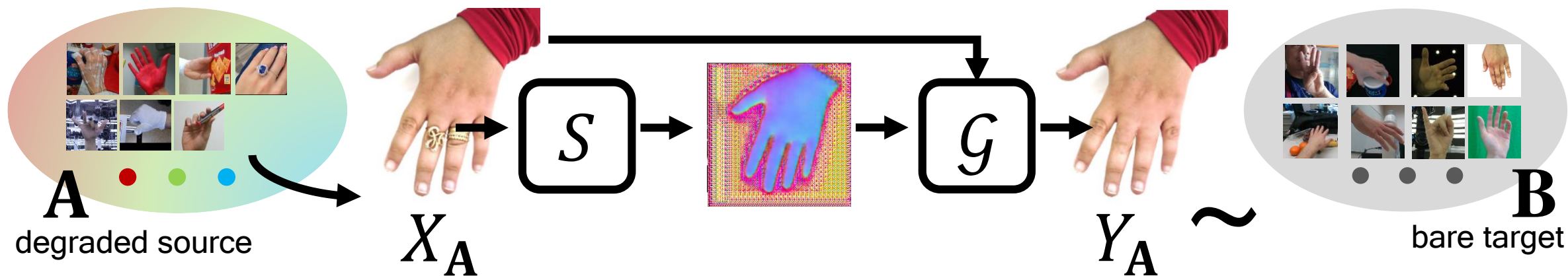
Gao et.al, NeurIPS 2022

- Pose & Appearance **Synchronous** Acquisition
 - Marker-less Procedure: Depends on the estimator **stability**
 - Synthetic Procedure: Suffers from the rendering **realness**

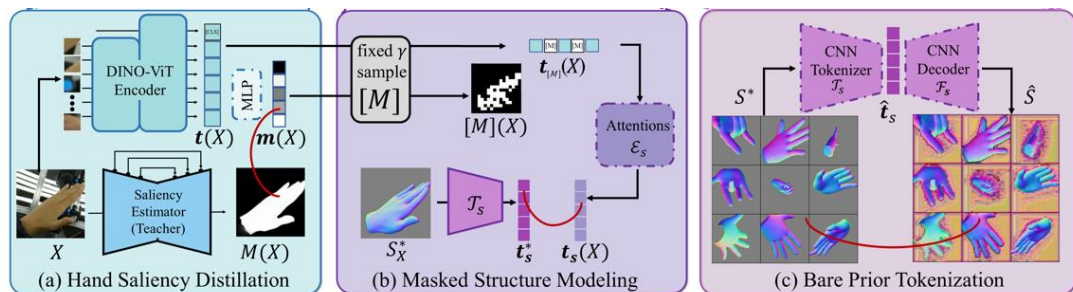
Insight



Overview

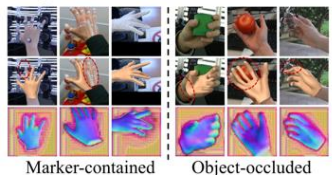


Structure Disentanglement

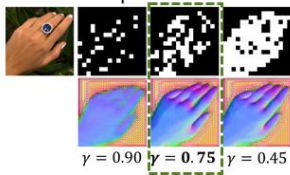


Sketcher: $X \rightarrow t(X) \rightarrow t_{[M]}(X) \rightarrow t_s(X) \rightarrow S(X)$

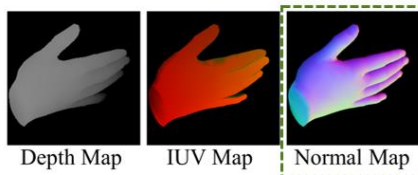
Degraded cases



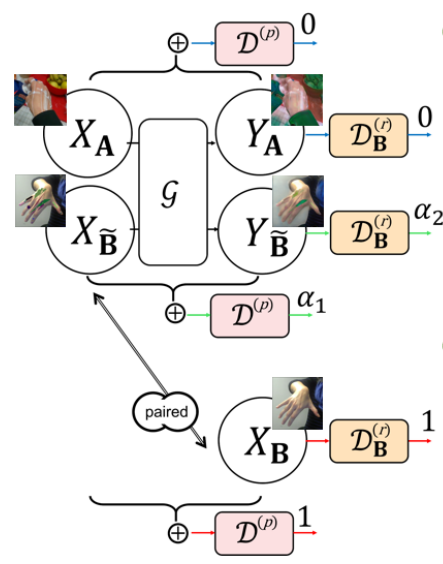
Ratio impact



Prior candidates



Dual Adversarial Discrimination



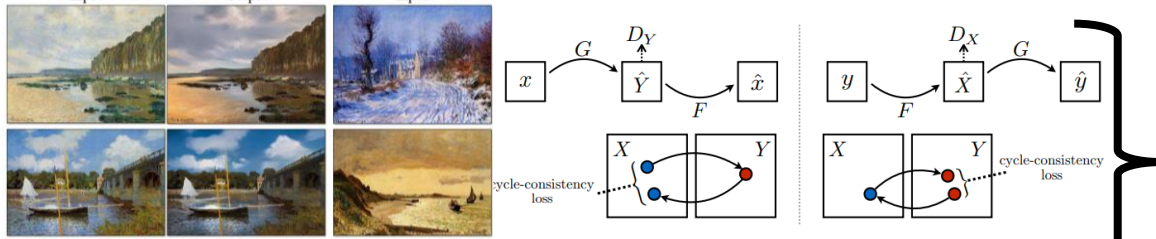
• Translation loss

$$L_G = \|(Y_A - X_A) \odot (1 - M[S(X_A)])\|_F + \|(Y_{\tilde{B}} - X_{\tilde{B}}) \odot (1 - M[S(X_{\tilde{B}})])\|_F + |\mathcal{D}_B^{(r)}(Y_A) - 1| + |\mathcal{D}_B^{(r)}(Y_{\tilde{B}}) - 1| + |\mathcal{D}^{(p)}(X_A \oplus Y_A) - 1| + |\mathcal{D}^{(p)}(X_{\tilde{B}} \oplus Y_{\tilde{B}}) - 1|$$

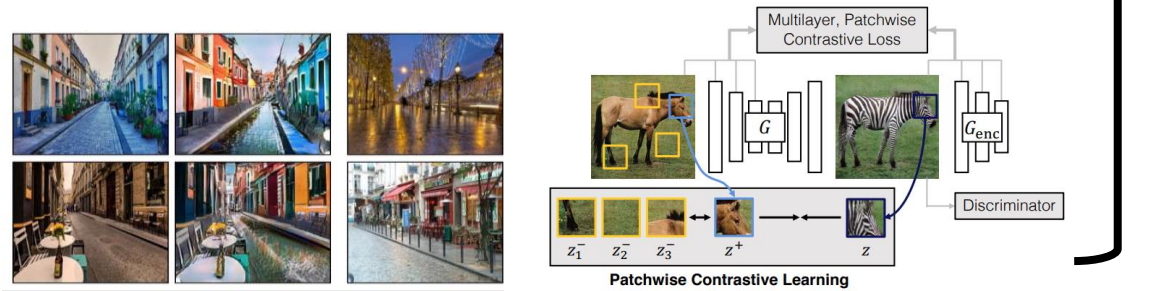
• Dual-discrimination loss

$$\begin{cases} L_D^{(r)} = |\mathcal{D}_B^{(r)}(Y_A) - 0| + |\mathcal{D}_B^{(r)}(Y_{\tilde{B}}) - \alpha_2| + |\mathcal{D}_B^{(r)}(X_B) - 1| \\ L_D^{(p)} = |\mathcal{D}^{(p)}(X_A \oplus Y_A) - 0| + |\mathcal{D}^{(p)}(X_{\tilde{B}} \oplus Y_{\tilde{B}}) - \alpha_1| + |\mathcal{D}^{(p)}(X_{\tilde{B}} \oplus X_B) - 1| \end{cases}$$

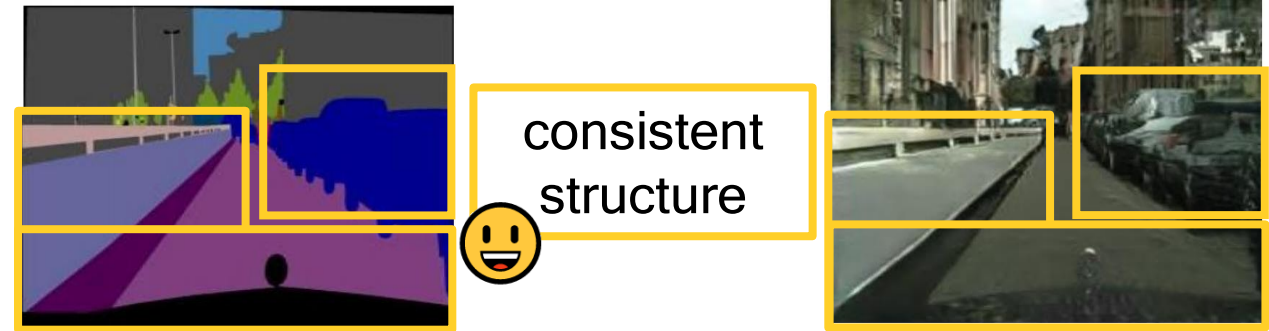
Challenges



CycleGAN, Zhu et.al. ICCV 2017



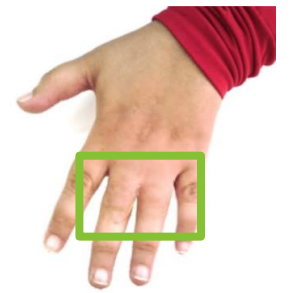
CUT, Park et.al. ECCV 2020



consistent structure

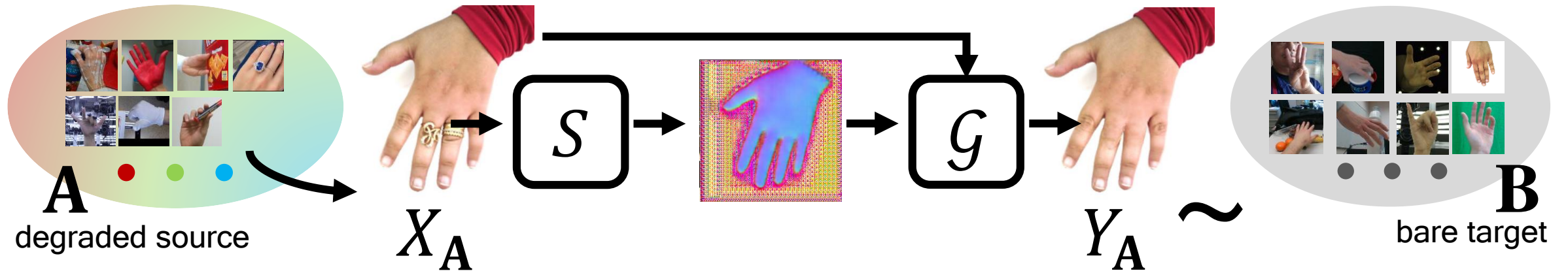


inconsistent structure



- **Data** : Degraded-bare hand image pairs are **almost impossible to obtain**.
- **Problem** : Background, foreground hand, and foreground markers require **different processing options**.
- **Baseline** : Existing unsupervised translators cannot handle the problem with **partial structure inconsistency**.

Contributions



- A semi-supervised framework that makes degraded images in marker-based MoCap regain bare appearance.
- A powerful ViT sketcher that disentangles bare hand structure without parametric model dependencies.
- An adversarial scheme that promotes the degraded-to-bare appearance wrapping effectively.

Step 1. Structure Disentanglement



Degraded Cases We Concerned :

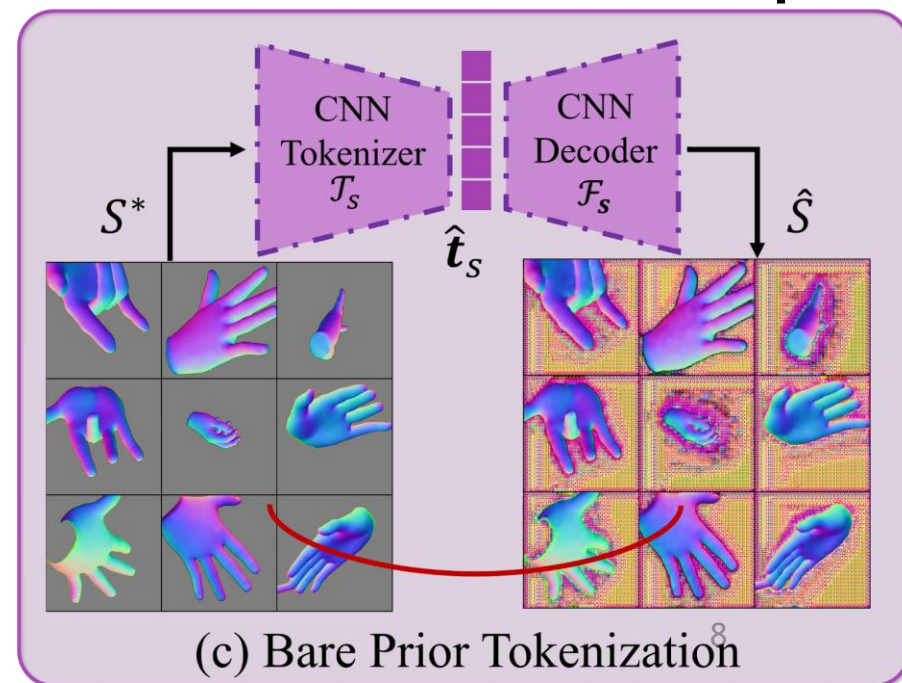
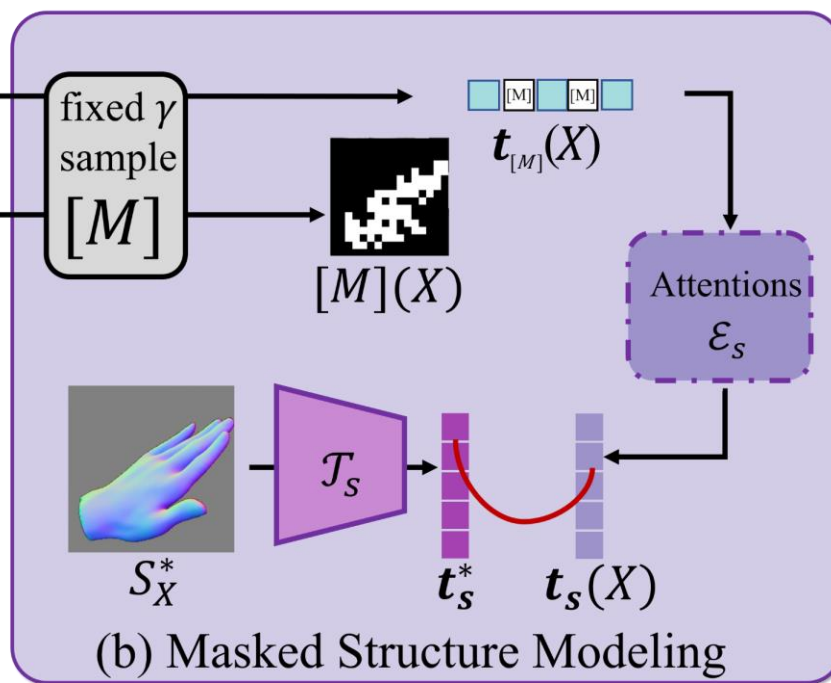
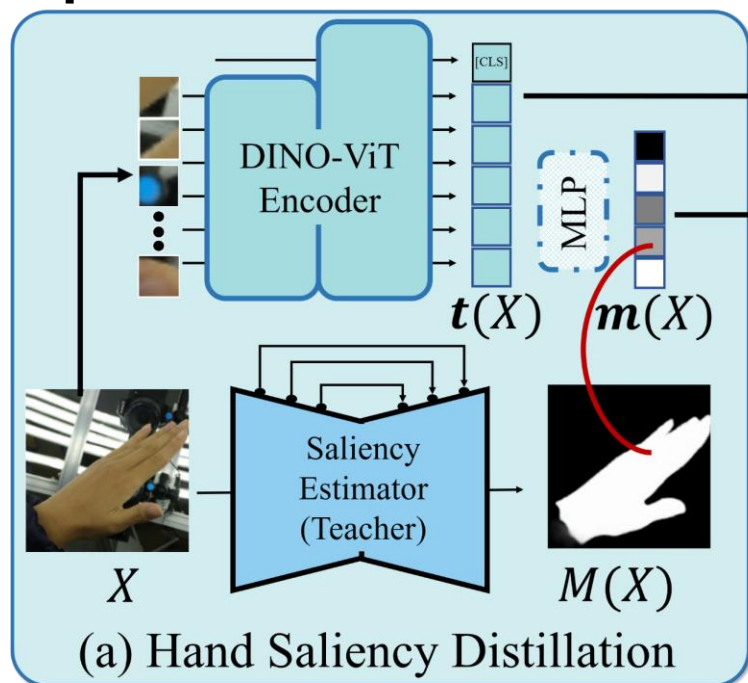
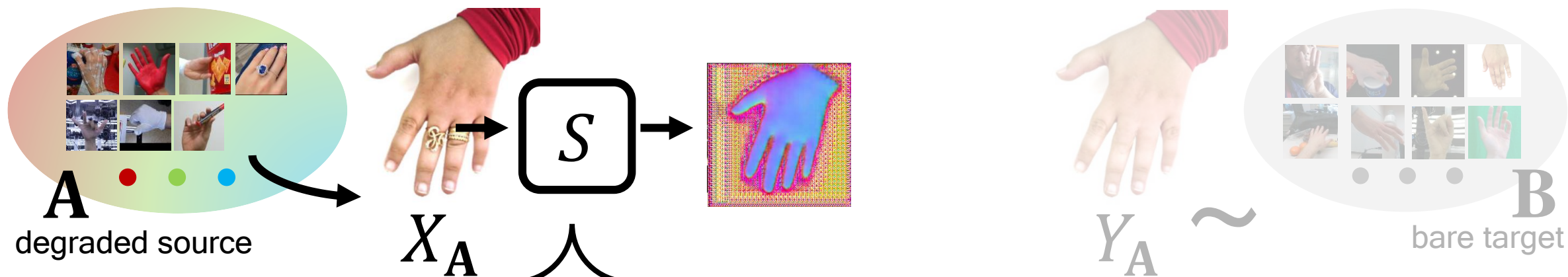


Marker-contained

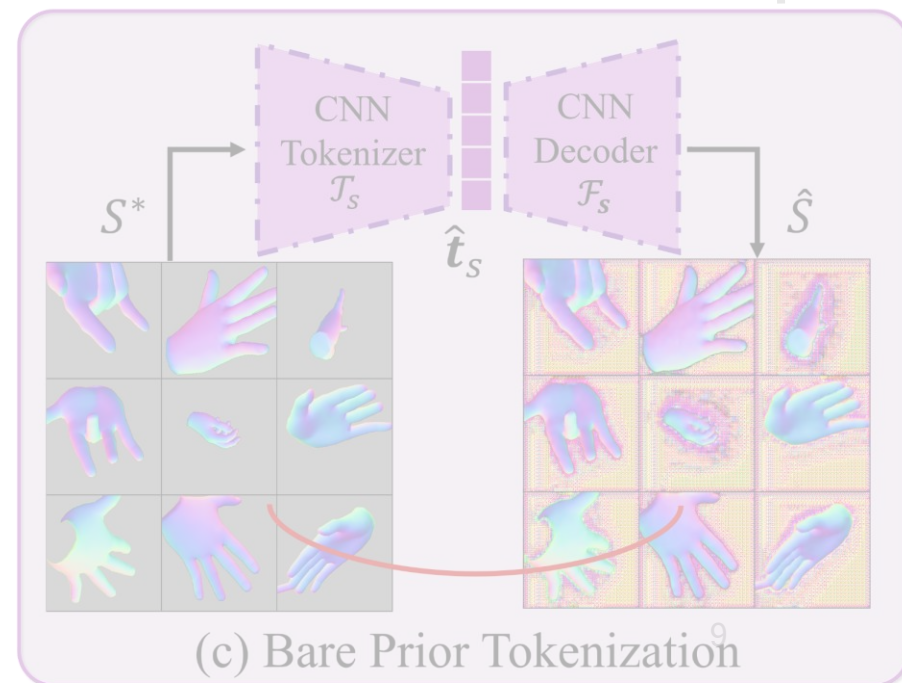
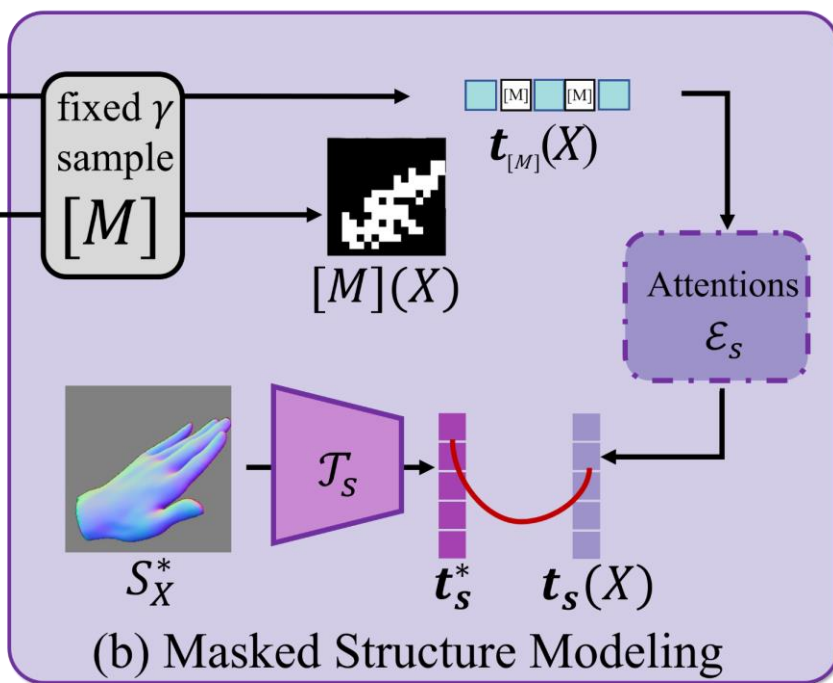
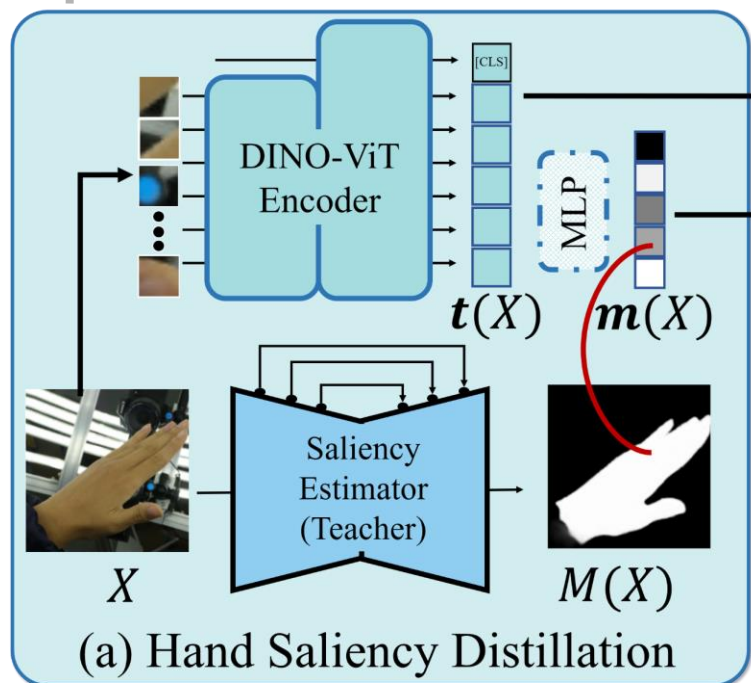
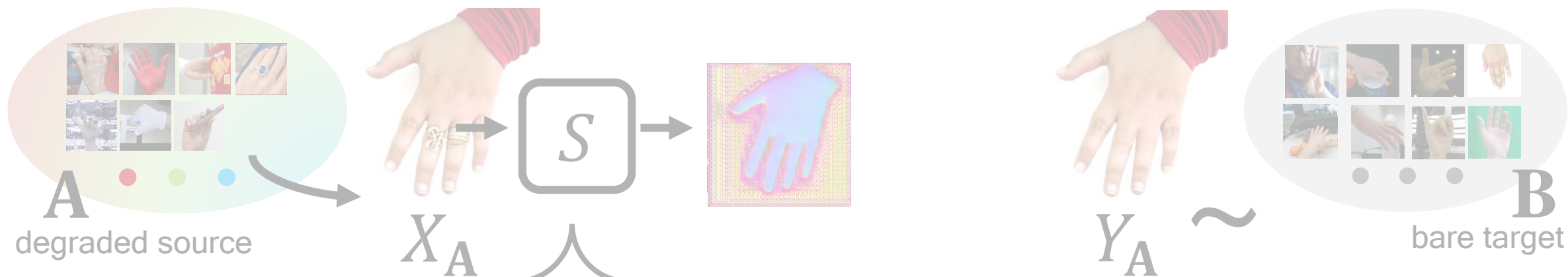


Object-occluded

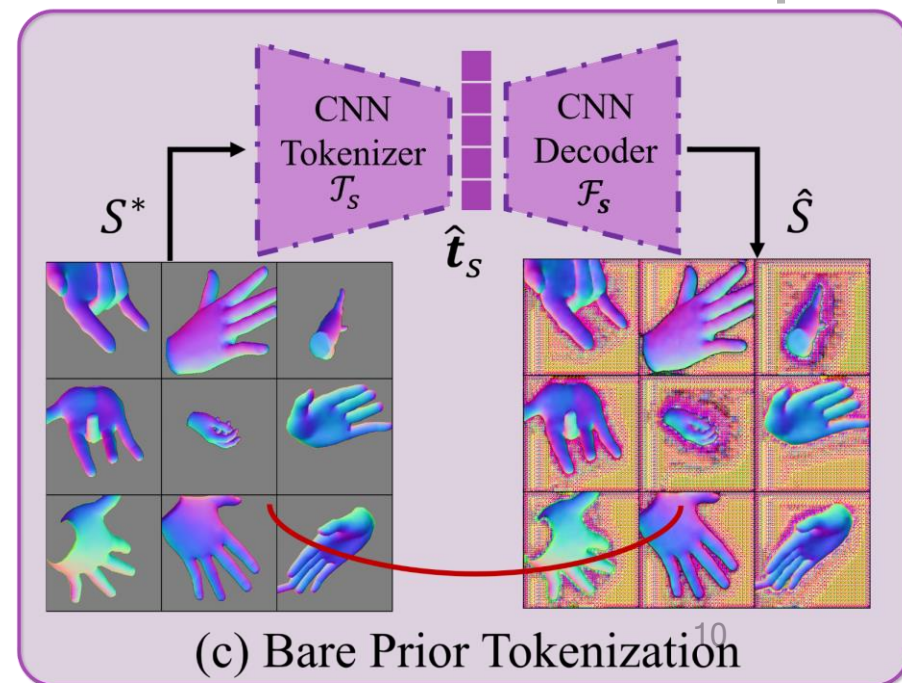
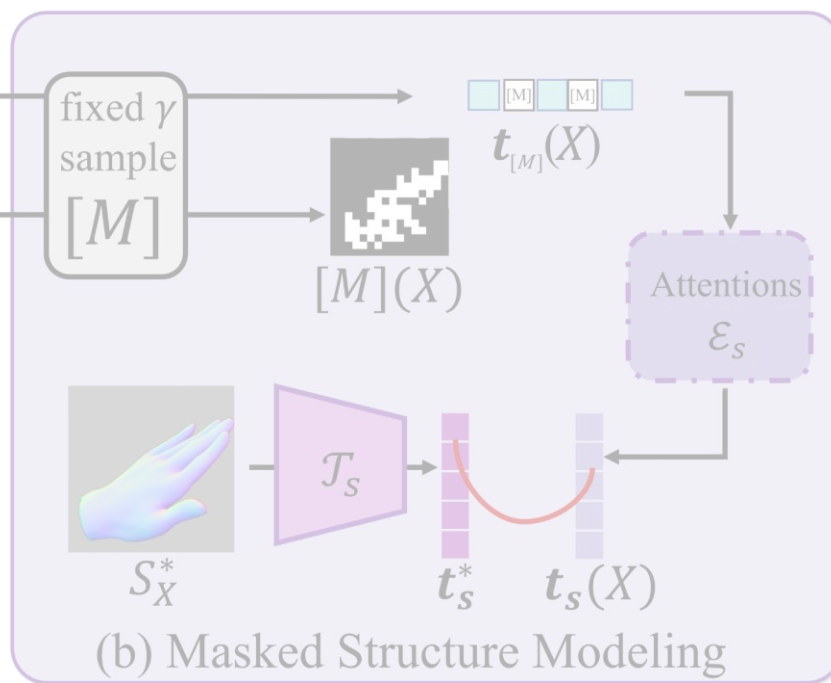
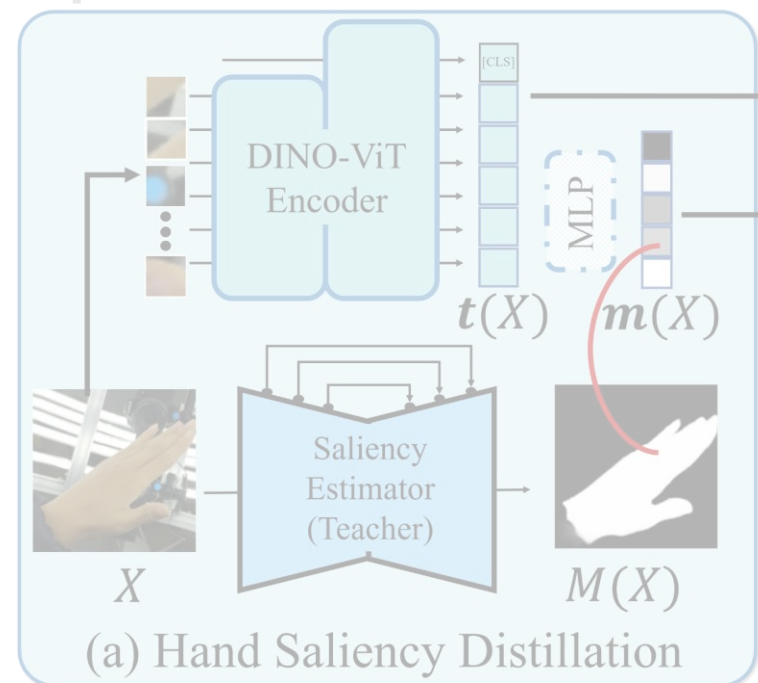
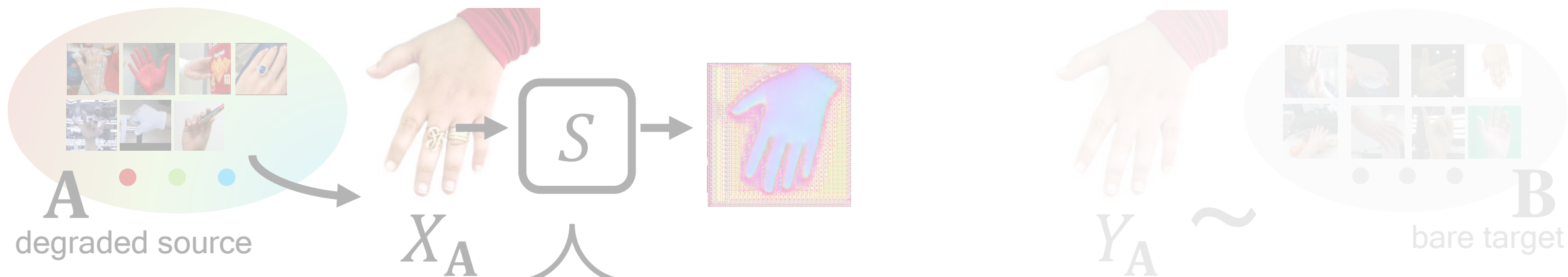
Step 1. Structure Disentanglement



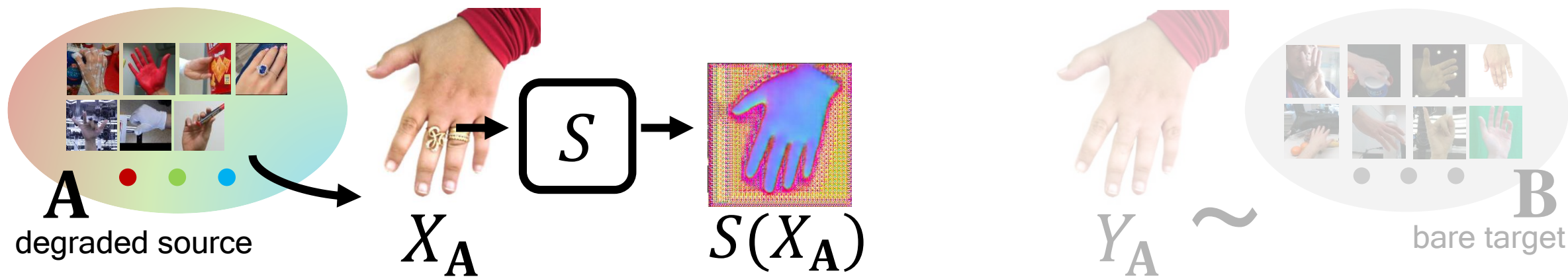
Step 1. Structure Disentanglement



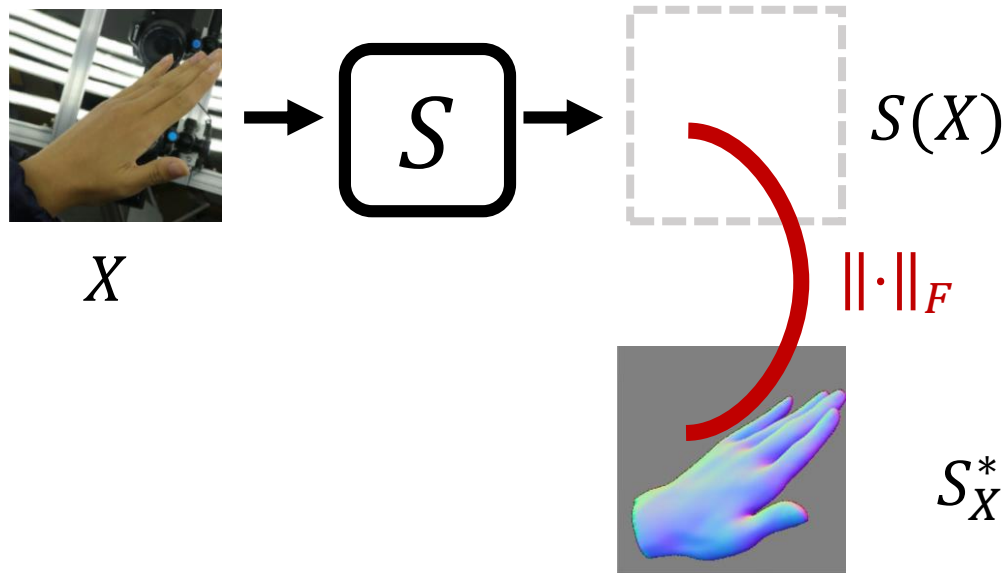
Step 1. Structure Disentanglement



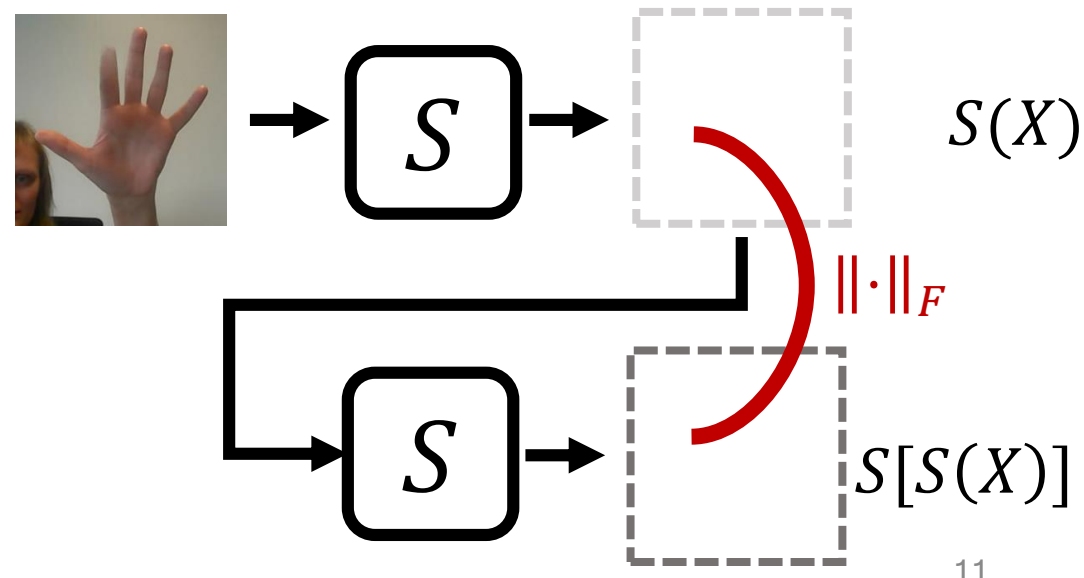
Step 1. Structure Disentanglement



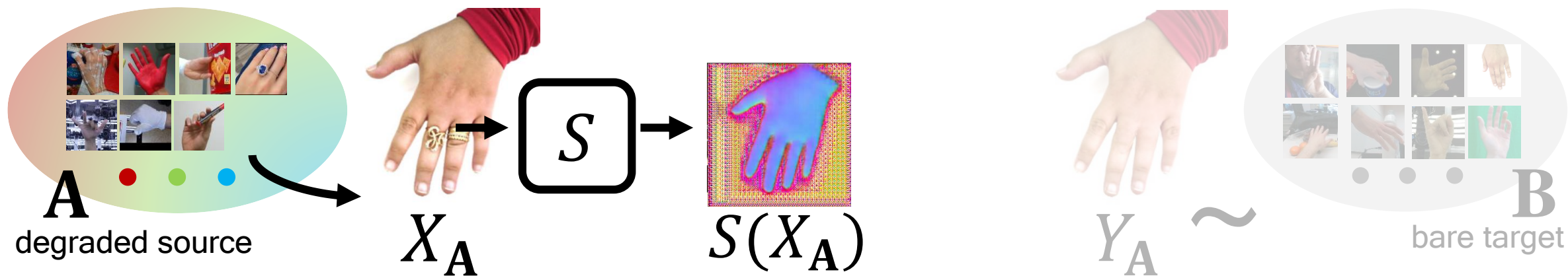
Train with **labeled** data $\{X, S_X^*\}$



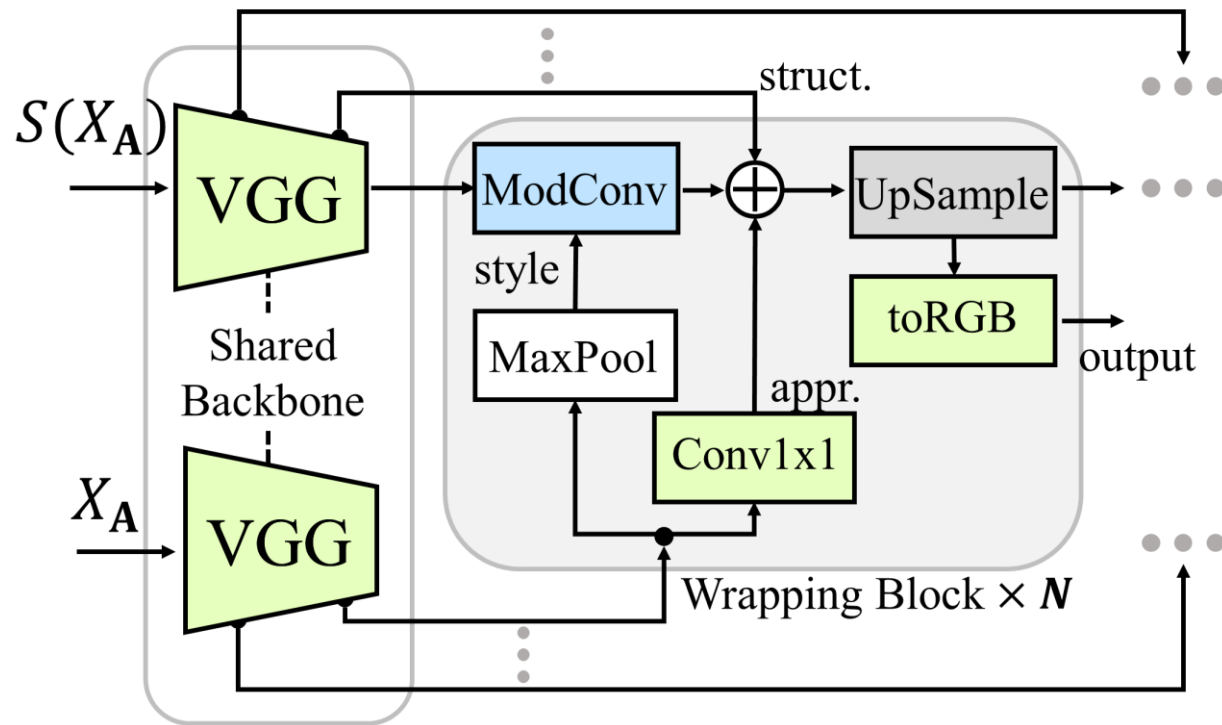
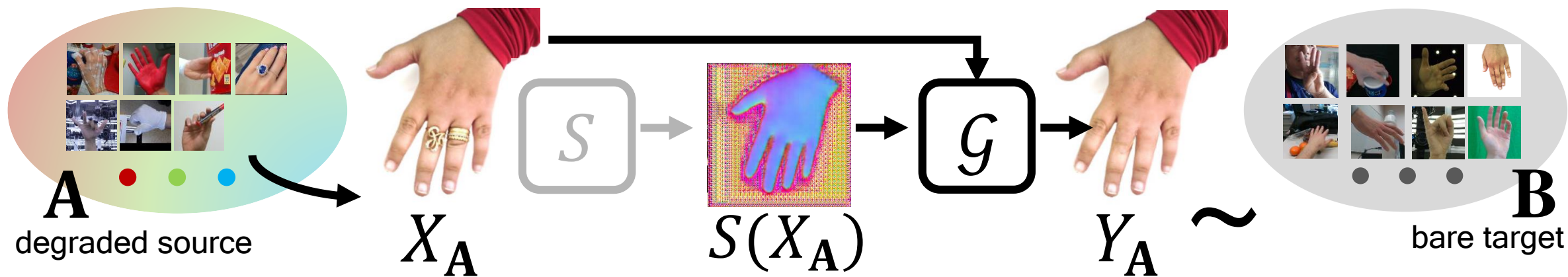
Train with **unlabeled** data $\{X\}$



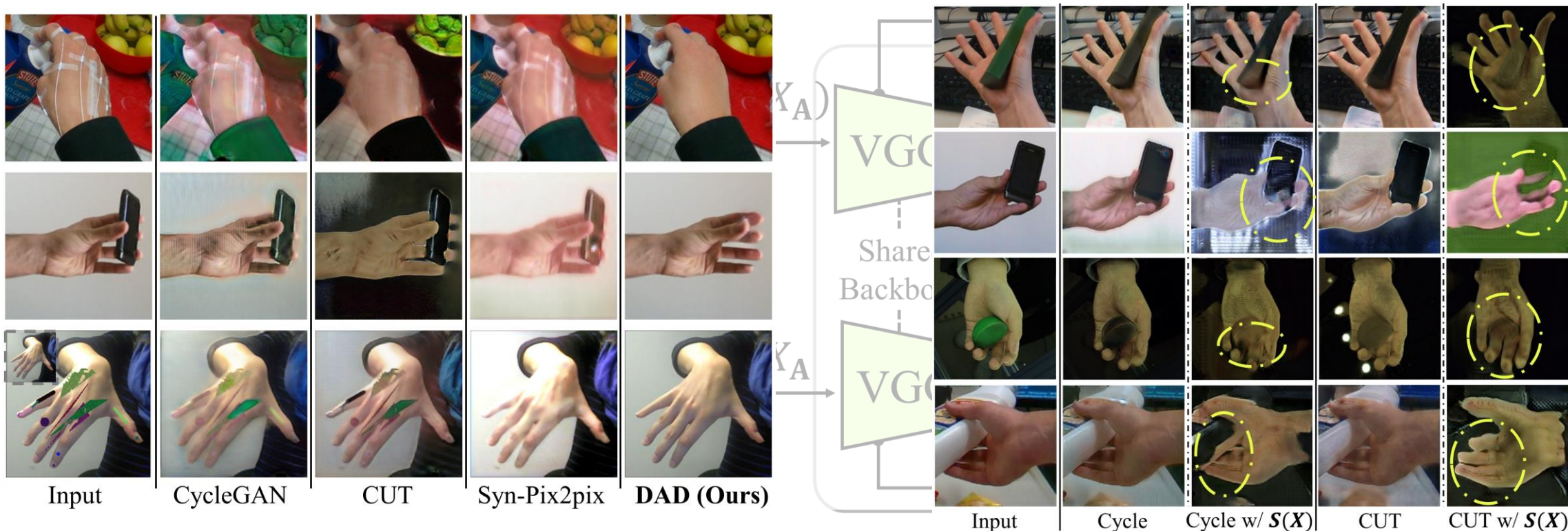
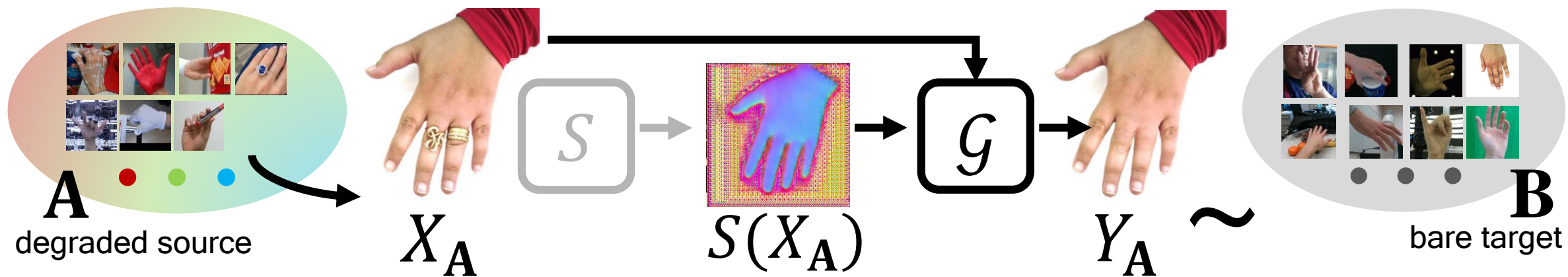
Step 1. Structure Disentanglement



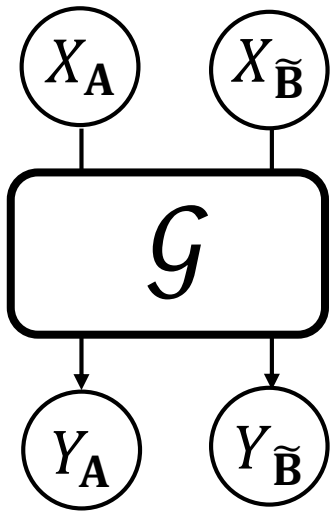
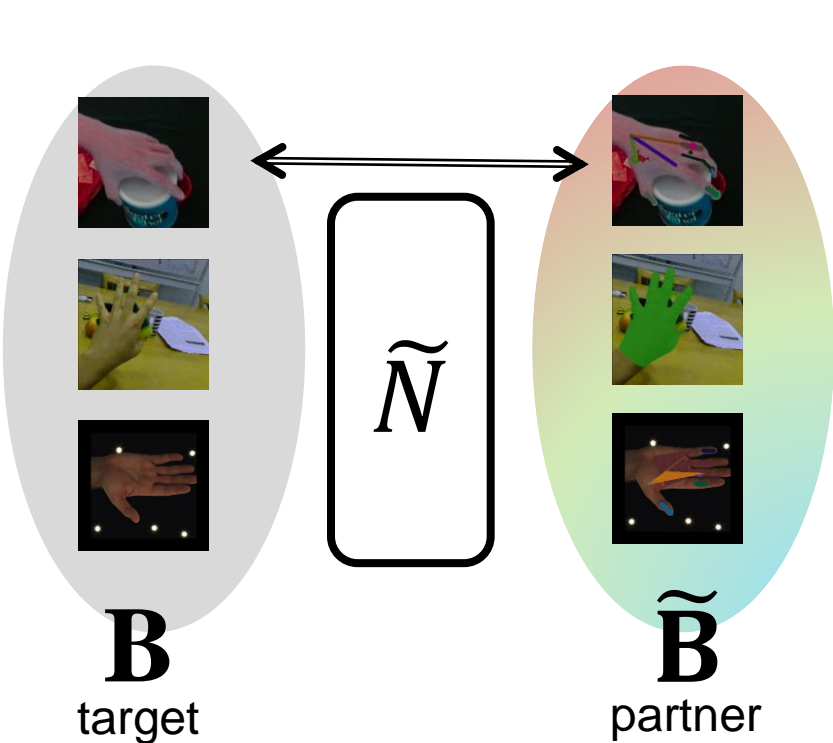
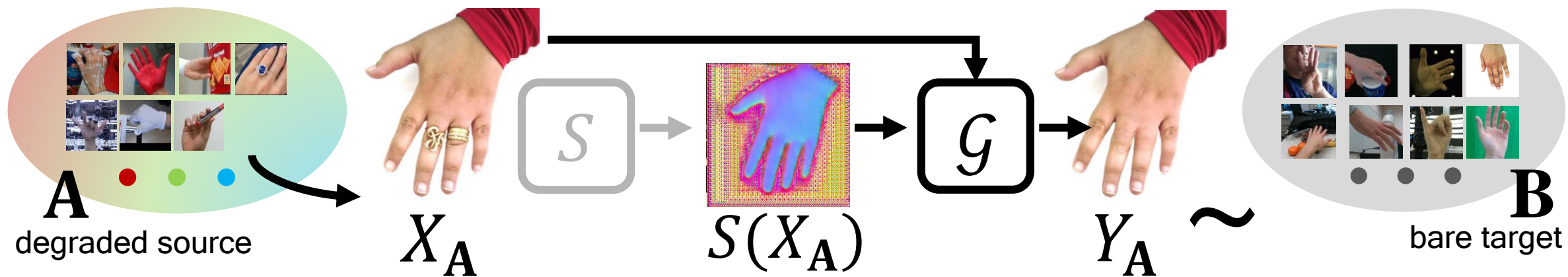
Step 2. Appearance Wrapping



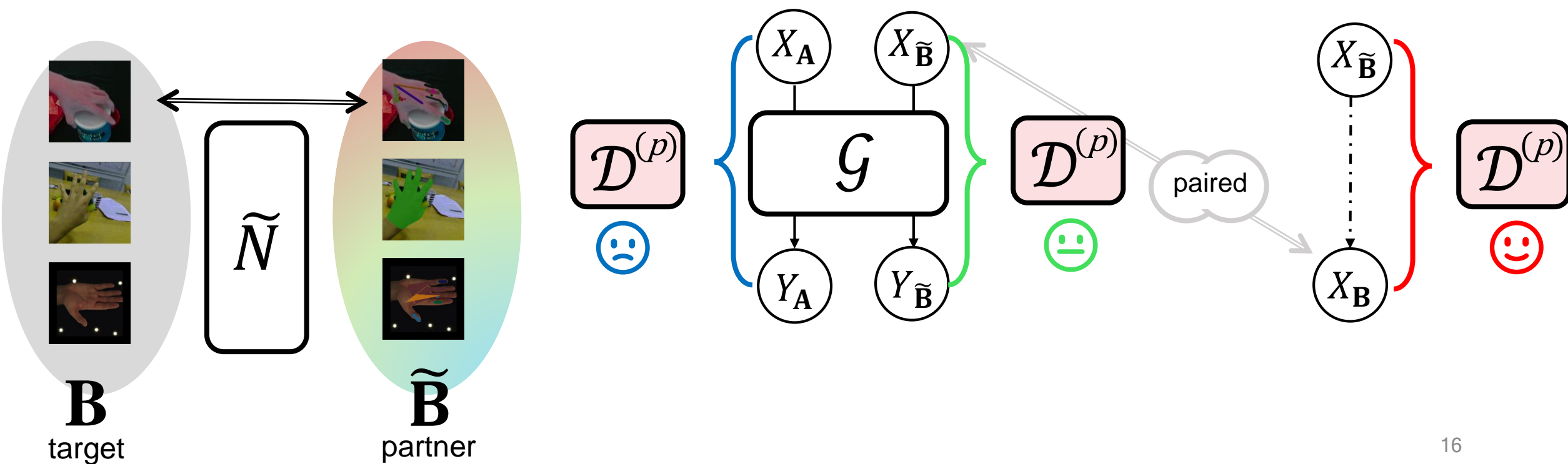
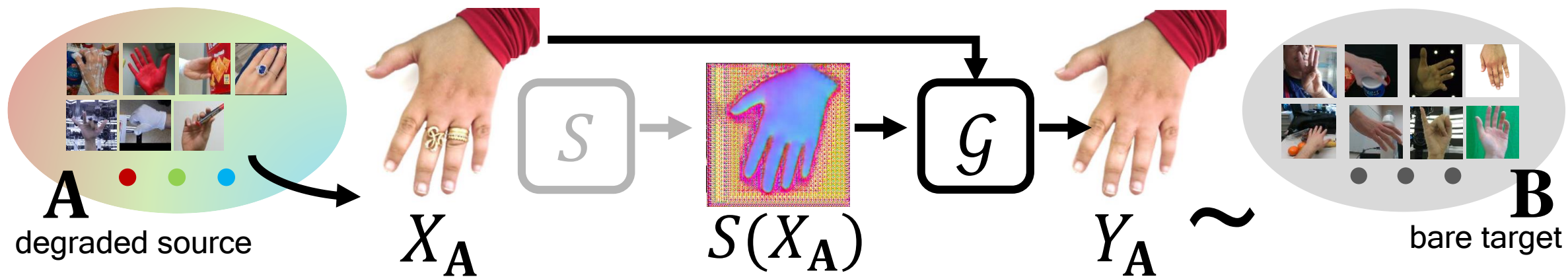
Step 2. Appearance Wrapping



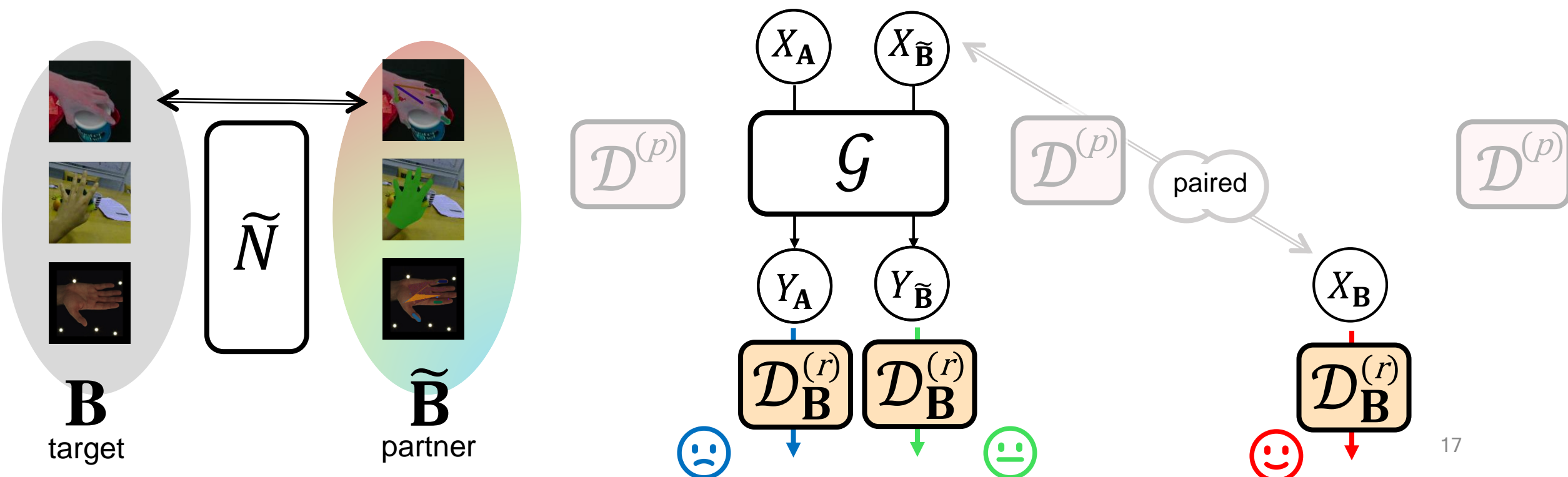
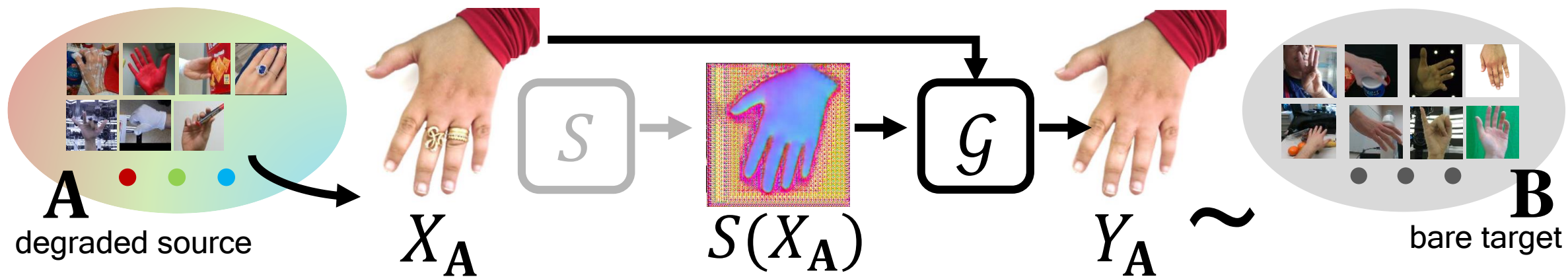
Step 2. Appearance Wrapping



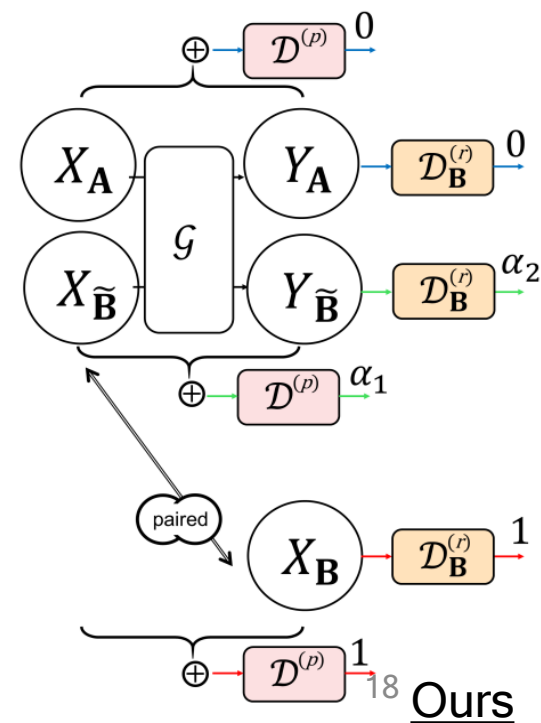
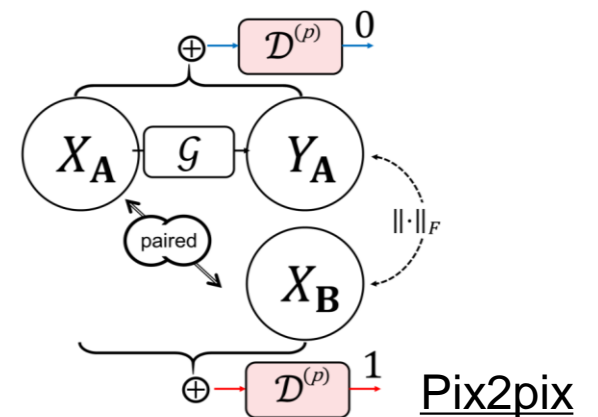
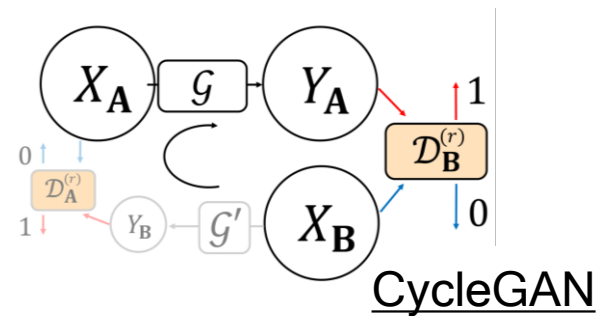
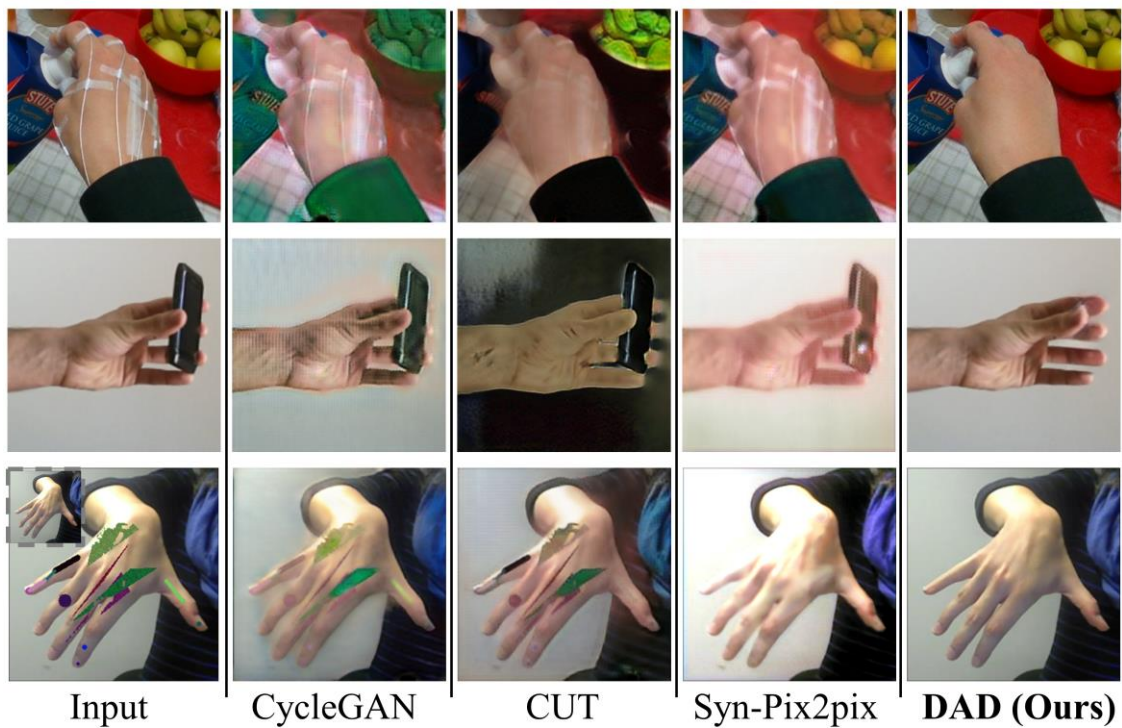
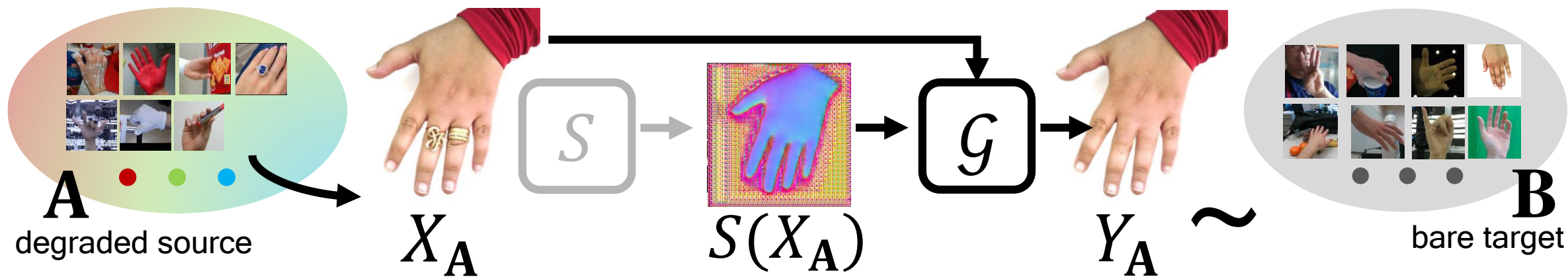
Step 2. Appearance Wrapping



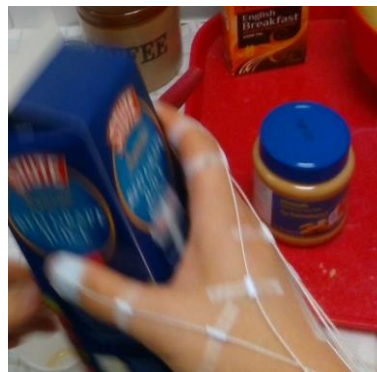
Step 2. Appearance Wrapping



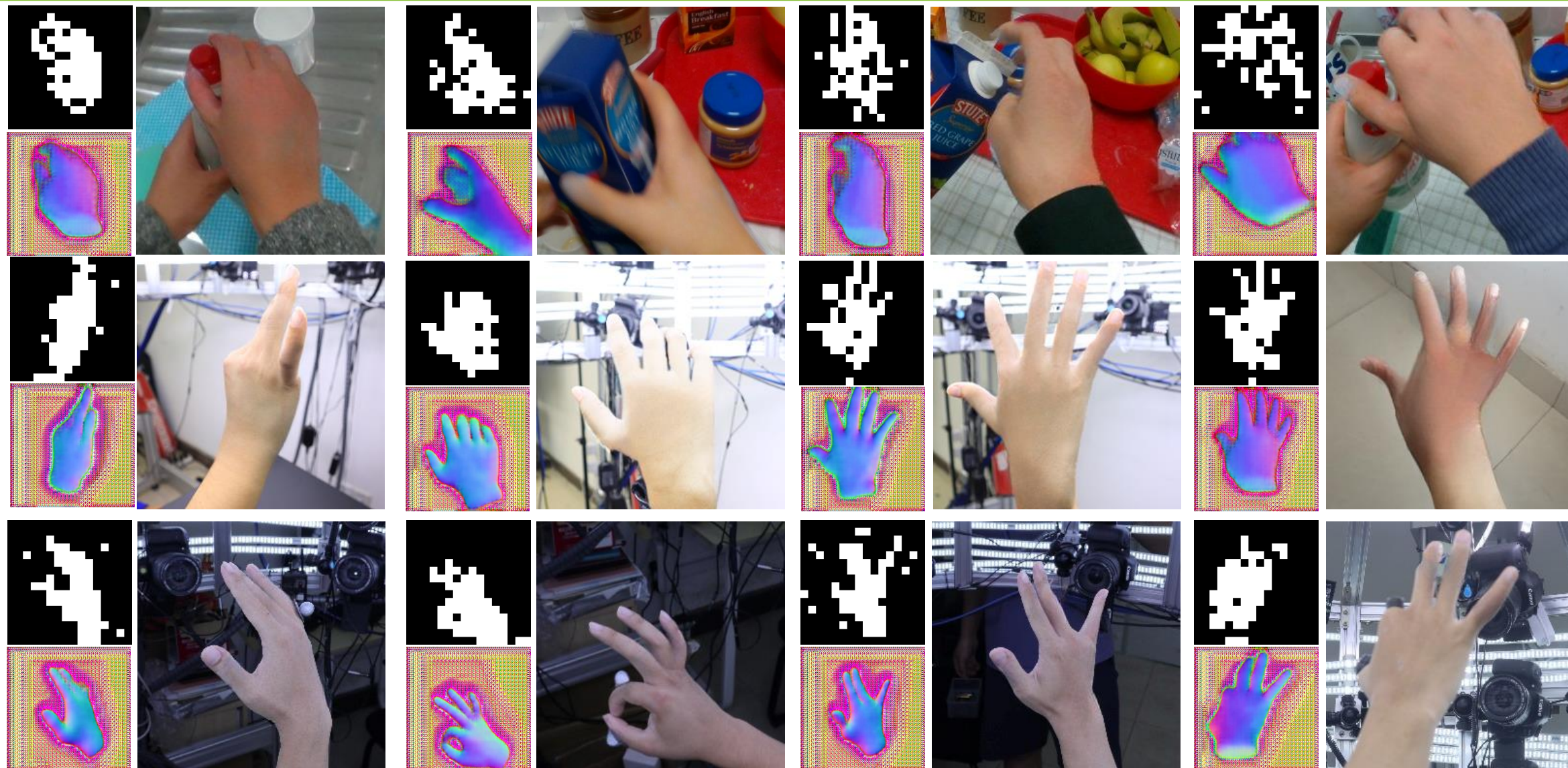
Step 2. Appearance Wrapping



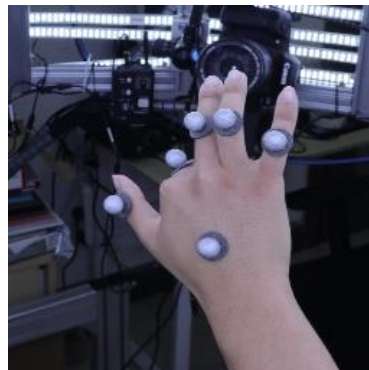
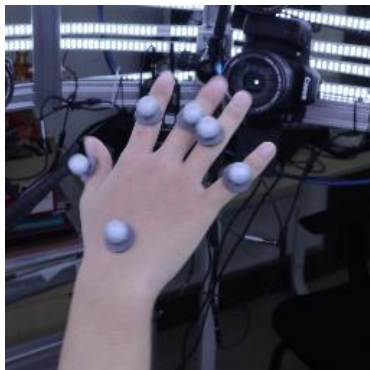
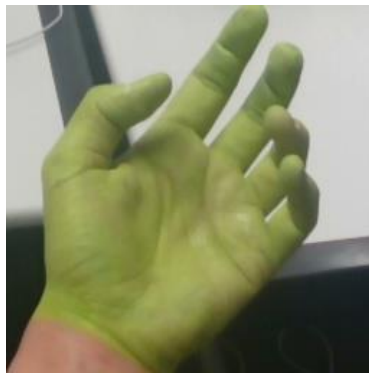
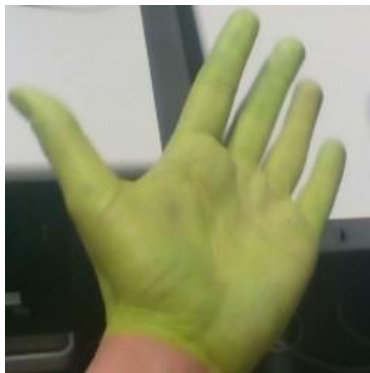
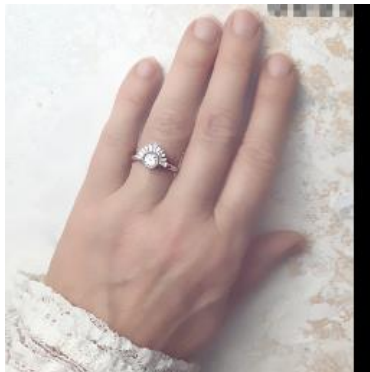
Qualitative Results on Translation $A_1 \rightarrow B$



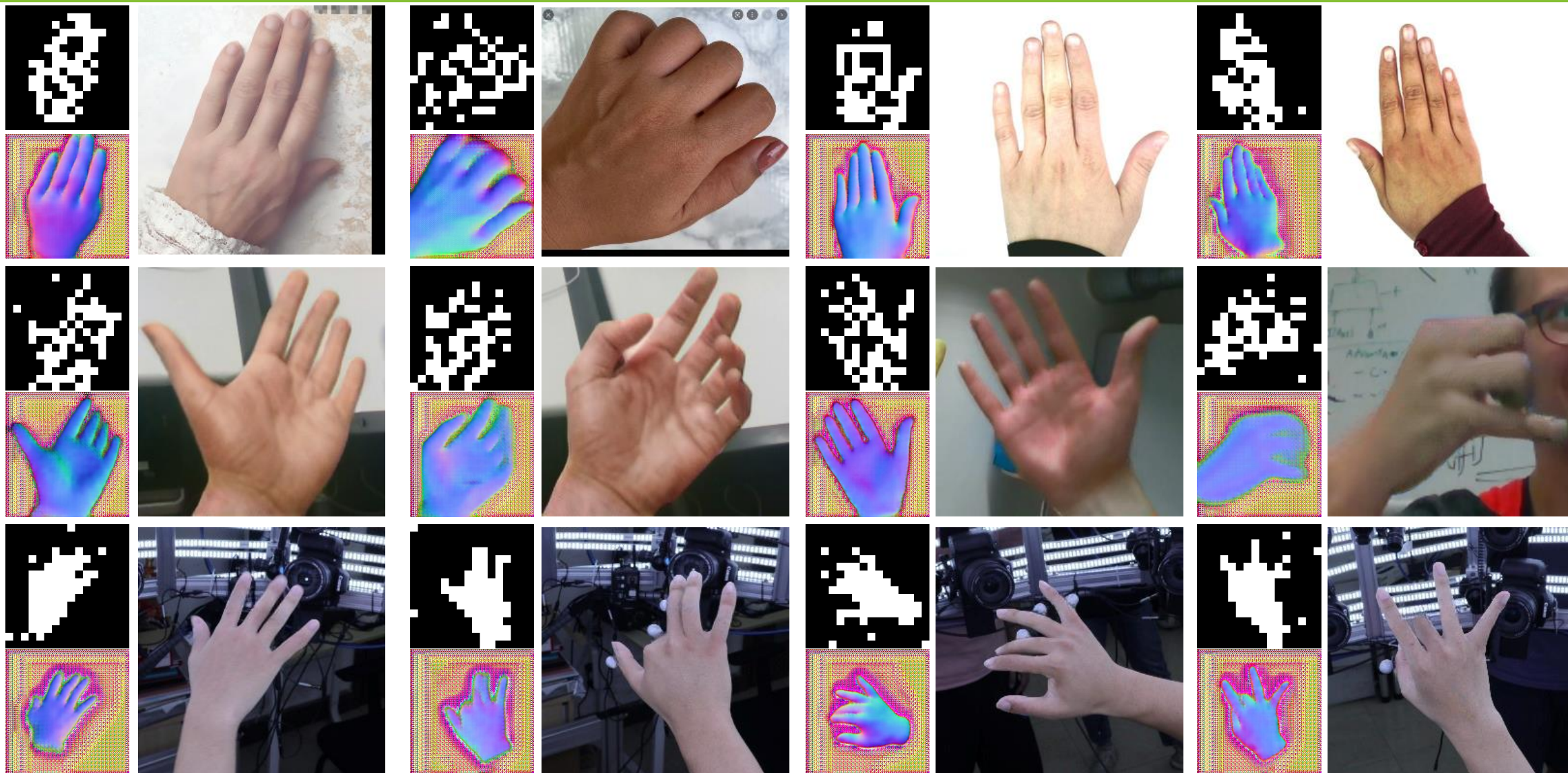
Qualitative Results on Translation $A_1 \rightarrow B$



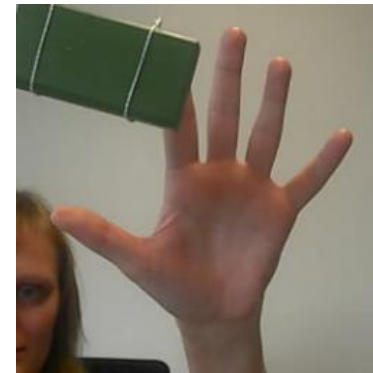
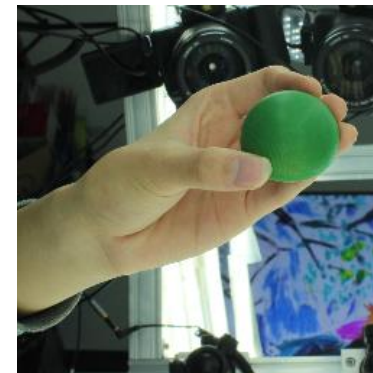
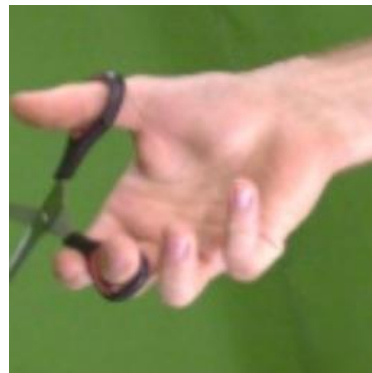
Qualitative Results on Translation $A_1 \rightarrow B$



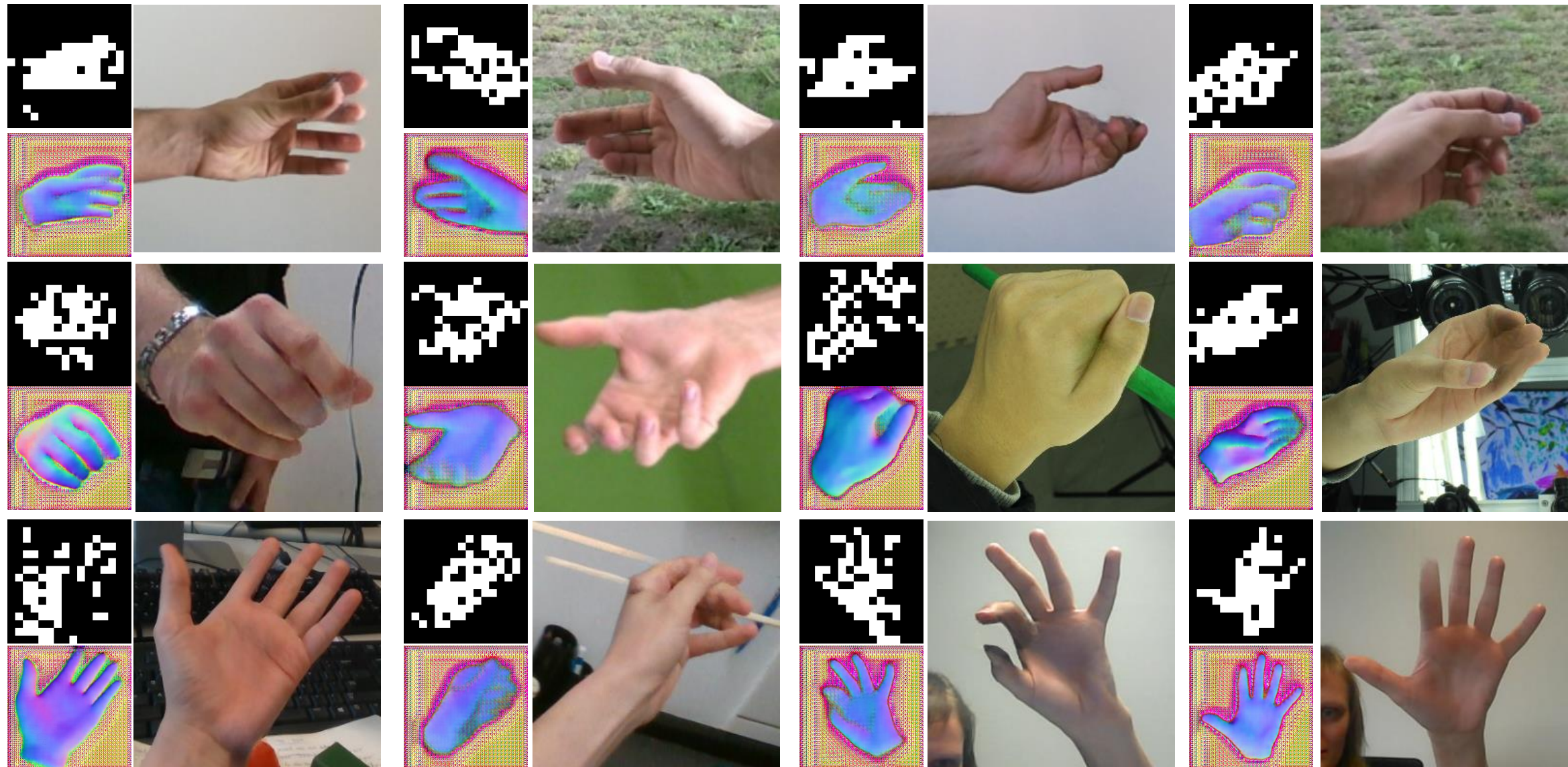
Qualitative Results on Translation $A_1 \rightarrow B$



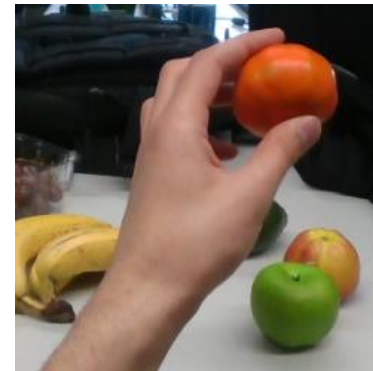
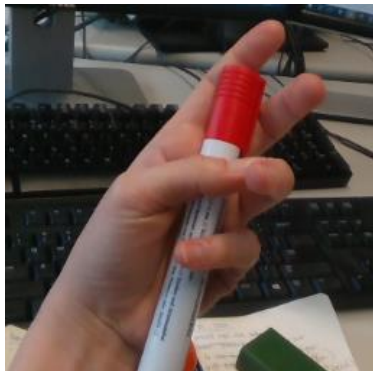
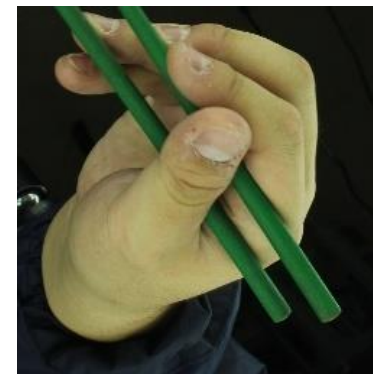
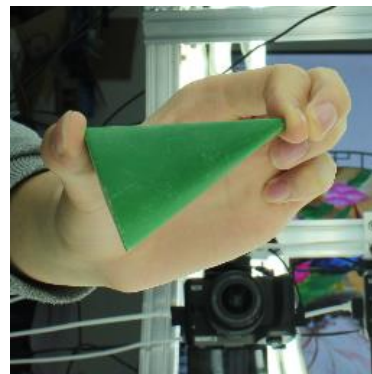
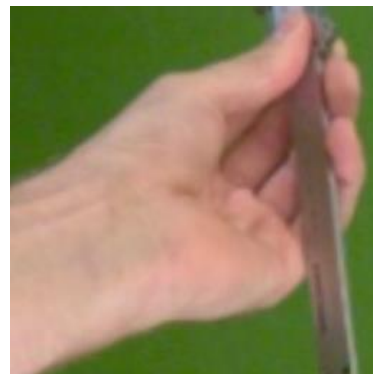
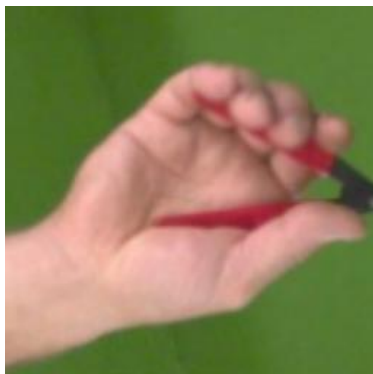
Qualitative Results on Translation $A_2 \rightarrow B$



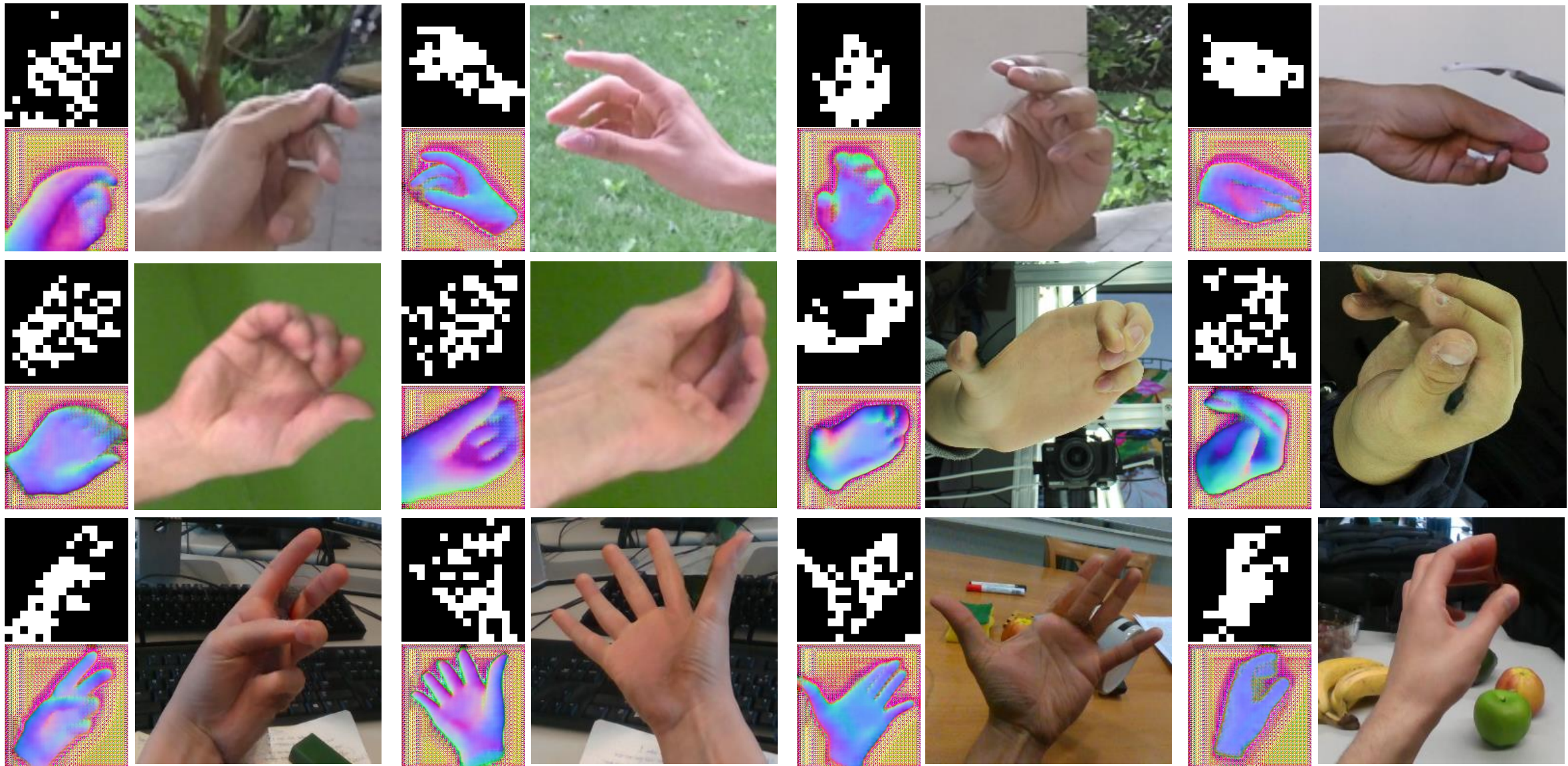
Qualitative Results on Translation $A_2 \rightarrow B$



Qualitative Results on Translation $A_2 \rightarrow B$



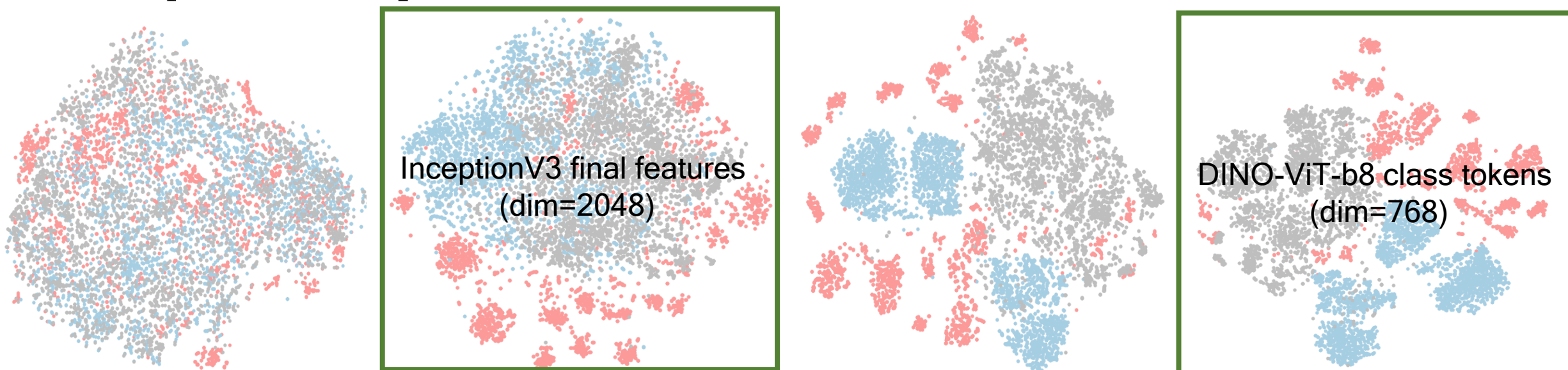
Qualitative Results on Translation $A_2 \rightarrow B$



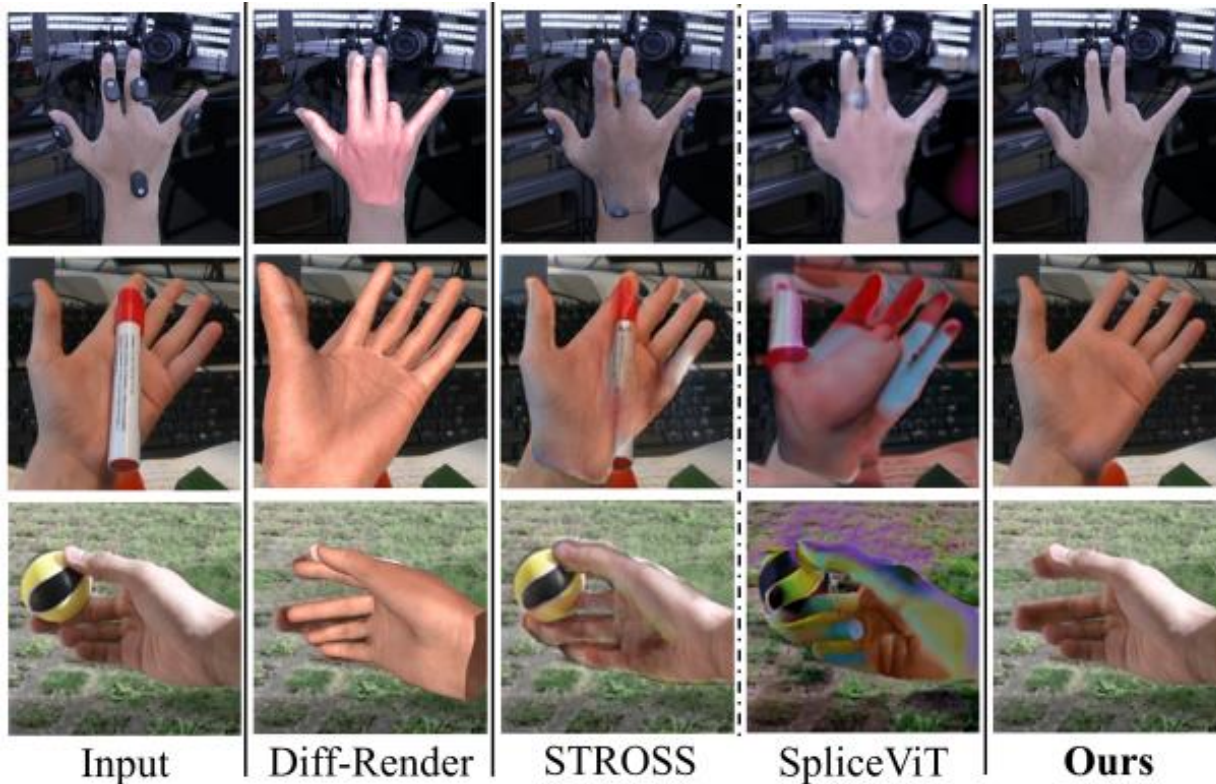
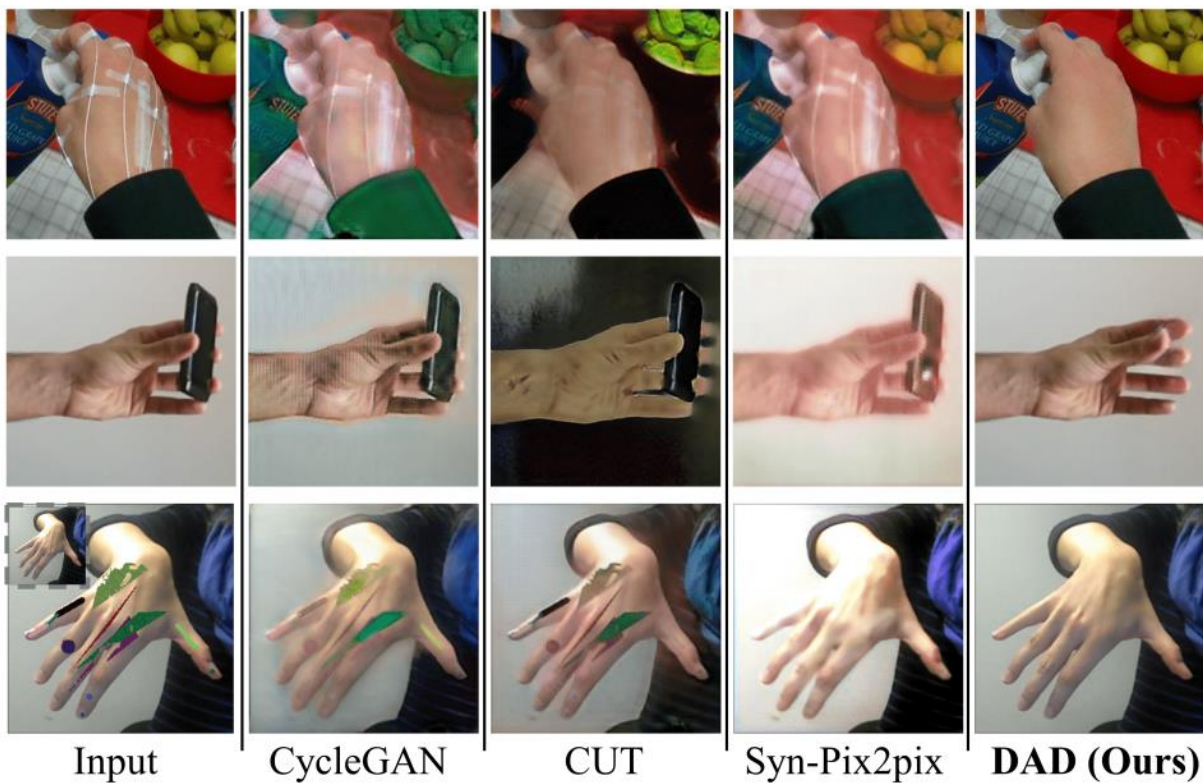
Quantitative Comparisons

| Tasks | $A_1 \rightarrow B$ | | | | $A_2 \rightarrow B$ | | | |
|----------------|---------------------|---------------------------|--------------------|----------------------|---------------------|---------------------------|--------------------|----------------------|
| Metrics | FID _i ↓ | KID _i (*100) ↓ | FID _v ↓ | KID _v ↓ | FID _i ↓ | KID _i (*100) ↓ | FID _v ↓ | KID _v ↓ |
| CycleGAN [91] | 76.39 | 4.46 ± 0.176 | 1266.17 | 32.02 ± 0.994 | 65.12 | 4.44 ± 0.196 | 1021.14 | 27.13 ± 0.948 |
| GANerated [52] | 76.97 | 4.72 ± 0.153 | 1220.53 | 31.72 ± 0.907 | 68.50 | 5.05 ± 0.199 | 985.82 | 26.75 ± 0.987 |
| H-GAN [56] | 93.02 | 6.53 ± 0.209 | 1488.32 | 37.85 ± 1.061 | 62.94 | 4.12 ± 0.197 | 876.94 | 24.61 ± 0.816 |
| UAG [5] | 87.80 | 5.90 ± 0.177 | 1375.55 | 35.67 ± 0.949 | 70.98 | 5.35 ± 0.210 | 1069.45 | 28.31 ± 0.926 |
| CUT [57] | 78.02 | 5.54 ± 0.192 | 1230.67 | 33.34 ± 1.015 | 58.88 | 3.73 ± 0.160 | 749.22 | 20.02 ± 0.826 |
| Ours | 60.37 | 3.45 ± 0.236 | 994.67 | 28.67 ± 0.916 | 41.53 | 3.37 ± 0.154 | 673.43 | 15.72 ± 1.209 |

● data in A_1 ● data in A_2 ● data in B



Qualitative Comparisons



Conclusion

- This work pioneers a **semi-supervised image-to-image translation** to recover the hand appearance that was originally degraded during the marker-based MoCap process.
- The **prior-based** sketcher can robustly disentangle the bare structure maps from degraded hand images.
- The DAD is more **efficient** than unsupervised schemes, and more **generalizable** than a supervised scheme trained with synthetic degradation.

