

A. Supplementary Material

A.1. Weighting Factors in the Reconstruction Loss of VAE Variants

Name	Value	
	<i>BaseVAE / StructVAE</i>	<i>MultiObjVAE / FactorVAE</i>
w_{T-A}^{pos}	14.5	5
w_{RH}^{pos}	14.5	10
w_{LH}^{pos}	14.5	12
w_F^{pos}	14.5	4
w_M^{pos}	0	8
w_{T-A}^{vel}	0	3
w_F^{vel}	0	6
w_M^{vel}	0	10

Table 1. Weights used in $\mathcal{L}_{\text{recon}}$ depending on the VAE variant (Equations 2 and 3).

A.2. Qualitative Comparison of Reconstructed and Generated Sign Pose Sequences

Figure 1 presents: 1) reconstructed poses from a ground-truth sign pose sequence of the PHOENIX14T dataset using *MultiObjVAE* and *FactorVAE*, and 2) the generated sequences from the corresponding input sentence¹ using the latent diffusion model trained on each VAE latent space. Despite visually similar VAE reconstructions for both variants (notably with correct hands and arms motions), we observe meaningful differences between generated sequences.

¹”Und nun die Wettervorhersage für morgen, Donnerstag, den sechszwanzigsten November.” meaning ”And now the weather forecast for tomorrow, Thursday, the twenty-sixth of November.”

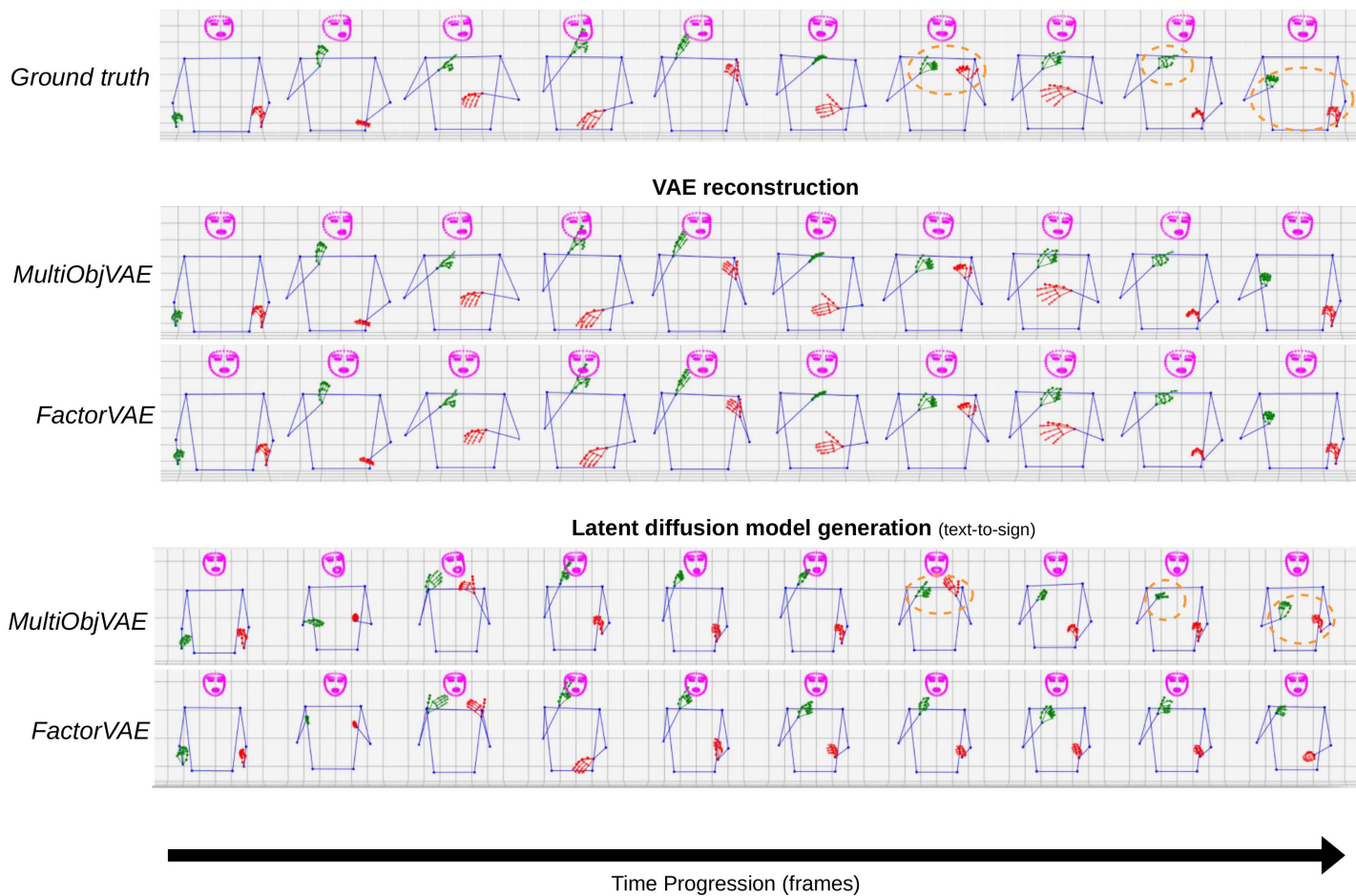


Figure 1. Qualitative comparison of sign pose sequences. Top: ground-truth poses and VAE reconstructions produced by *MultiObjVAE* and *FactorVAE*. Bottom: text-to-sign sequences generated by latent diffusion models trained on the corresponding latent spaces. Orange dotted circles indicate hand configurations correctly generated in the *MultiObjVAE* latent space but not in the *FactorVAE* case. Sequence ID in the PHOENIX14T dataset: 25November_2009_Wednesday_tagesschau-7666.