

A. Theoretical Analysis of the Trade-off Between Pretrained and One-Point Velocities for Oracle MeanFlow Learning

Proposition A.1. For any $\lambda \in [0, 1]$, consider the combination of the one-point estimator and the pre-trained teacher velocity

$$\mathbf{w}_\lambda := (1 - \lambda)\hat{\mathbf{v}} + \lambda\mathbf{v}_\phi.$$

Plugging \mathbf{w}_λ into the target $\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w})$ (with $\mathbf{w} = \mathbf{w}_\lambda$) yields the corresponding loss $\mathcal{L}_{\text{MF}}(\theta; \mathbf{w}_\lambda)$. Consider the following three residuals: the one-point velocity residual, the teacher-oracle velocity residual, and the oracle bias.

$$\begin{aligned}\delta\hat{\mathbf{v}}_t &:= \hat{\mathbf{v}}(\mathbf{z}_t, t) - \mathbf{v}(\mathbf{z}_t, t), \\ \delta\mathbf{v}_t^\phi &:= \mathbf{v}_\phi(\mathbf{z}_t, t) - \mathbf{v}(\mathbf{z}_t, t), \\ \delta\mathbf{h}(\mathbf{z}_t, t, s) &:= \mathbf{h}_{\theta^-}(\mathbf{z}_t, t, s) - \mathbf{h}(\mathbf{z}_t, t, s).\end{aligned}$$

Then the MF loss admits the following decomposition:

$$\mathcal{L}_{\text{MF}}(\theta; \mathbf{w}_\lambda) = \mathbb{E}_{t, \mathbf{z}_t} \left\| \mathbf{h}_{\theta^-} - (\mathbf{h} + \mathbf{B} + \lambda \mathbf{A}_{\theta^-} \delta\mathbf{v}_t^\phi) \right\|^2 + (1 - \lambda)^2 \mathbb{E} \left\| \mathbf{A}_{\theta^-} \delta\hat{\mathbf{v}}_t \right\|^2,$$

where $\mathbf{B}(\mathbf{z}_t, t, s) := (t - s) \left(\partial_t \delta\mathbf{h} + (\nabla_{\mathbf{x}} \delta\mathbf{h}) \mathbf{v} \right)$ and $\mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s) := \mathbf{I} - (t - s) \nabla_{\mathbf{x}} \mathbf{h}_{\theta^-}$.

Proof. Throughout the proof we fix s and suppress the dependence on (\mathbf{z}_t, t, s) whenever there is no ambiguity. All expectations are taken over (t, \mathbf{z}_t) and when $\hat{\mathbf{v}}$ is random (e.g. a Monte-Carlo estimator).

We first derive the affine form of the teacher target around the oracle velocity \mathbf{v} . Recall from the MF construction that for any proxy velocity \mathbf{w} , the target associated with \mathbf{h}_{θ^-} is

$$\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}) = \mathbf{w} - (t - s) \left((\nabla_{\mathbf{x}} \mathbf{h}_{\theta^-}) \mathbf{w} + \partial_t \mathbf{h}_{\theta^-} \right),$$

so for two velocities \mathbf{w} and \mathbf{v} we have

$$\begin{aligned}\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}) - \mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{v}) &= \left[\mathbf{w} - (t - s) (\nabla_{\mathbf{x}} \mathbf{h}_{\theta^-}) \mathbf{w} \right] - \left[\mathbf{v} - (t - s) (\nabla_{\mathbf{x}} \mathbf{h}_{\theta^-}) \mathbf{v} \right] \\ &= (\mathbf{I} - (t - s) \nabla_{\mathbf{x}} \mathbf{h}_{\theta^-}) (\mathbf{w} - \mathbf{v}).\end{aligned}$$

By the definition of \mathbf{A}_{θ^-} in the proposition, $\mathbf{A}_{\theta^-} := \mathbf{I} - (t - s) \nabla_{\mathbf{x}} \mathbf{h}_{\theta^-}$, so this can be written as

$$\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}) = \mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{v}) + \mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s) (\mathbf{w} - \mathbf{v}(\mathbf{z}_t, t)).$$

On the other hand, by the PDE relation between the oracle flow map \mathbf{h} and the oracle velocity \mathbf{v} , and by the definition $\delta\mathbf{h} = \mathbf{h}_{\theta^-} - \mathbf{h}$ and $\mathbf{B}(\mathbf{z}_t, t, s) := (t - s) \left(\partial_t \delta\mathbf{h} + (\nabla_{\mathbf{x}} \delta\mathbf{h}) \mathbf{v} \right)$, we have (see the derivation in the main text)

$$\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{v}) = \mathbf{h}(\mathbf{z}_t, t, s) + \mathbf{B}(\mathbf{z}_t, t, s).$$

Combining the last two displays gives, for every proxy velocity \mathbf{w} ,

$$\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}) = \mathbf{h}(\mathbf{z}_t, t, s) + \mathbf{B}(\mathbf{z}_t, t, s) + \mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s) (\mathbf{w} - \mathbf{v}(\mathbf{z}_t, t)), \quad (5)$$

which is exactly the affine reparametrization of the target around the oracle velocity \mathbf{v} .

We now prove the claimed decomposition of the MF loss. By the definitions of the residual velocities,

$$\hat{\mathbf{v}} = \mathbf{v} + \delta\hat{\mathbf{v}}_t, \quad \mathbf{v}_\phi = \mathbf{v} + \delta\mathbf{v}_t^\phi,$$

the mixed velocity satisfies, for any $\lambda \in [0, 1]$,

$$\mathbf{w}_\lambda = (1 - \lambda)\hat{\mathbf{v}} + \lambda\mathbf{v}_\phi = (1 - \lambda)(\mathbf{v} + \delta\hat{\mathbf{v}}_t) + \lambda(\mathbf{v} + \delta\mathbf{v}_t^\phi) = \mathbf{v} + (1 - \lambda)\delta\hat{\mathbf{v}}_t + \lambda\delta\mathbf{v}_t^\phi. \quad (6)$$

Substituting $\mathbf{w} = \mathbf{w}_\lambda$ and Equation (6) into Equation (5) yields

$$\begin{aligned}\mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}_\lambda) &= \mathbf{h} + \mathbf{B} + \mathbf{A}_{\theta^-}(\mathbf{w}_\lambda - \mathbf{v}) \\ &= \mathbf{h} + \mathbf{B} + \mathbf{A}_{\theta^-}((1 - \lambda)\delta\hat{\mathbf{v}}_t + \lambda\delta\mathbf{v}_t^\phi) \\ &= \mathbf{h}(\mathbf{z}_t, t, s) + \mathbf{B}(\mathbf{z}_t, t, s) + \lambda\mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s)\delta\mathbf{v}_t^\phi + (1 - \lambda)\mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s)\delta\hat{\mathbf{v}}_t.\end{aligned}\quad (7)$$

For the remainder of the proof we fix (t, \mathbf{z}_t) and abbreviate

$$\mathbf{h}_\theta := \mathbf{h}_\theta(\mathbf{z}_t, t, s), \quad \mathbf{h} := \mathbf{h}(\mathbf{z}_t, t, s), \quad \mathbf{A}_{\theta^-} := \mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s), \quad \mathbf{B} := \mathbf{B}(\mathbf{z}_t, t, s).$$

By Equation (2), the MF loss at \mathbf{w}_λ is

$$\mathcal{L}_{\text{MF}}(\theta; \mathbf{w}_\lambda) = \mathbb{E} \left\| \mathbf{h}_\theta(\mathbf{z}_t, t, s) - \mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}_\lambda) \right\|_2^2.$$

Using Equation (7), we can write the pointwise residual as

$$\mathbf{h}_\theta - \mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}_\lambda) = \left[\mathbf{h}_\theta - (\mathbf{h} + \mathbf{B} + \lambda\mathbf{A}_{\theta^-}\delta\mathbf{v}_t^\phi) \right] - (1 - \lambda)\mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t.$$

Let

$$\mathbf{Y} := \mathbf{h}_\theta - (\mathbf{h} + \mathbf{B} + \lambda\mathbf{A}_{\theta^-}\delta\mathbf{v}_t^\phi),$$

so that

$$\mathbf{h}_\theta - \mathbf{h}_{\theta^-}^{\text{tgt}}(\mathbf{z}_t, t, s; \mathbf{w}_\lambda) = \mathbf{Y} - (1 - \lambda)\mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t.$$

Therefore

$$\mathcal{L}_{\text{MF}}(\theta; \mathbf{w}_\lambda) = \mathbb{E} \left\| \mathbf{Y} - (1 - \lambda)\mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t \right\|_2^2 = \mathbb{E} \|\mathbf{Y}\|_2^2 - 2(1 - \lambda)\mathbb{E} \langle \mathbf{Y}, \mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t \rangle + (1 - \lambda)^2 \mathbb{E} \|\mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t\|_2^2. \quad (8)$$

Since the one-point estimator is conditionally unbiased:

$$\mathbb{E}[\hat{\mathbf{v}}(\mathbf{z}_t, t) \mid \mathbf{z}_t] = \mathbf{v}(\mathbf{z}_t, t), \quad \text{so} \quad \mathbb{E}[\delta\hat{\mathbf{v}}_t \mid \mathbf{z}_t] = \mathbf{0}.$$

For fixed (t, \mathbf{z}_t) , the quantities $\mathbf{h}_\theta, \mathbf{h}, \mathbf{B}, \mathbf{v}, \mathbf{v}_\phi$, hence \mathbf{Y} and \mathbf{A}_{θ^-} , are deterministic, and the only randomness comes from the internal noise of $\hat{\mathbf{v}}$. Thus

$$\begin{aligned}\mathbb{E} \langle \mathbf{Y}, \mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t \rangle &= \mathbb{E}_{t, \mathbf{z}_t} \left[\mathbb{E}[\langle \mathbf{Y}, \mathbf{A}_{\theta^-}\delta\hat{\mathbf{v}}_t \rangle \mid \mathbf{z}_t] \right] \\ &= \mathbb{E}_{t, \mathbf{z}_t} \left[\langle \mathbf{Y}, \mathbf{A}_{\theta^-} \mathbb{E}[\delta\hat{\mathbf{v}}_t \mid \mathbf{z}_t] \rangle \right] \\ &= \mathbb{E}_{t, \mathbf{z}_t} \langle \mathbf{Y}, \mathbf{A}_{\theta^-} \cdot \mathbf{0} \rangle = \mathbf{0}.\end{aligned}$$

Hence the cross term in Equation (8) is zero, and we obtain

$$\mathcal{L}_{\text{MF}}(\theta; \mathbf{w}_\lambda) = \mathbb{E}_{t, \mathbf{z}_t} \left\| \mathbf{h}_\theta(\mathbf{z}_t, t, s) - (\mathbf{h}(\mathbf{z}_t, t, s) + \mathbf{B}(\mathbf{z}_t, t, s) + \lambda\mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s)\delta\mathbf{v}_t^\phi) \right\|_2^2 + (1 - \lambda)^2 \mathbb{E} \|\mathbf{A}_{\theta^-}(\mathbf{z}_t, t, s)\delta\hat{\mathbf{v}}_t\|_2^2,$$

which is exactly the claimed decomposition Equation (4). This completes the proof. \square

From Proposition 3.1, several implications follow. First, using the pre-trained velocity \mathbf{v}_ϕ improves stability: when $\lambda = 1$, the noisy one-point term vanishes, which reduces gradient variance and typically stabilizes and accelerates optimization.

Second, this comes at the cost of bias. The model no longer regresses to \mathbf{h} , but to $\mathbf{h} + \mathbf{B} + \lambda\mathbf{A}_{\theta^-}\delta\mathbf{v}_s^\phi$. When \mathbf{v}_ϕ is accurate (small $\delta\mathbf{v}_s^\phi$) and the time step is small so that $\mathbf{A}_{\theta^-} \approx \mathbf{I}$, this additional bias is small. In the small-step regime, the shift is approximately $\lambda\delta\mathbf{v}_s^\phi$, so the practical target is close to the oracle one whenever $\|\delta\mathbf{v}_s^\phi\|$ is small. However, under domain shift, $\delta\mathbf{v}_s^\phi$ can be non-negligible and the learned \mathbf{h}_θ will partially compensate for this mismatch.

Third, the mixing weight λ offers a natural tuning mechanism. A schedule with λ gradually increasing to 1 can start from a regime with low bias and low variance: initially, the one-point label maintains an unbiased target while the pre-trained signal reduces variance; later, larger λ leverages the smoothness of \mathbf{v}_ϕ for stable convergence. In practice, we use a simple two-stage schedule: $\lambda = 0$ in the distillation stage and $\lambda = 1$ in the (optional) bootstrapping stage.

Finally, the coupling to the self-teacher remains simple. Both \mathbf{B} and \mathbf{A}_{θ^-} are treated as stop-gradient with respect to θ^- , so updates in θ still solve a least-squares fit to a shifted target. As the self-teacher improves over training, $\delta\mathbf{h}$ shrinks and \mathbf{B} decreases, further reducing the induced bias.

B. More Hyperparameter Settings

MF-RAE’s Flow Matching Pre-Training Setup. We directly use the flow matching model checkpoint provided by RAE [48], where the model on ImageNet 256 is trained for 800 epochs, and that on ImageNet 512 is trained for 400 epochs.

MF-RAE’s CMT Setup. The CMT stage’s learning rate is also $1e-4$, but the EMA β is set to 0.999 for faster convergence, thanks to CMT’s stability and fixed target.

MF-RAE’s CMT Setup on ImageNet 256. CMT generates eight trajectories from the teacher DiT^{DH}-XL without guidance during each iteration. We use the Euler ODE solver with 16 steps (FID=2.3). For each trajectory, it generates $16 \times 15/2 = 120$ pairs, but we randomly select 64 pairs for optimization to fit on the H100 GPU. The total batch size is thus $64 \times 8 = 512$ pairs. We conduct this training on 8 H100 GPUs for 27K iterations.

MF-RAE’s CMT Setup on ImageNet 512. CMT generates two trajectories from the teacher DiT^{DH}-XL without guidance during each iteration. We use the Euler ODE solver with 16 steps (FID=1.66 without any guidance). For each trajectory, we randomly select 8 pairs for optimization to fit on the H100 GPU. The total batch size is thus $8 \times 2 = 16$ pairs. We accumulate gradients for eight rounds to enlarge the batch size. We conduct this training on 2 H100 GPUs for 27K iterations.

C. Generated Samples



Figure 2. ImageNet 256 MF-RAE 1-step samples on class 437: beacon, lighthouse, beacon light, pharos.

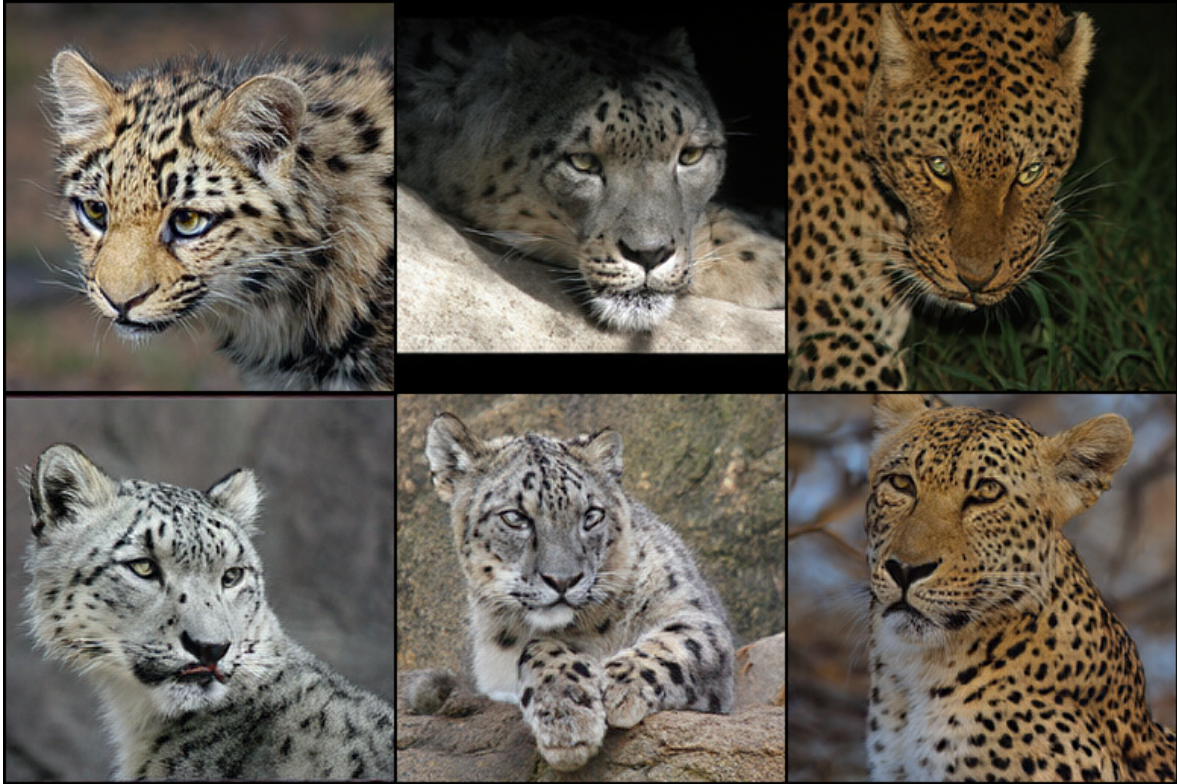


Figure 3. ImageNet 256 MF-RAE 1-step samples on classes 288 and 290: leopard and snow leopard.



Figure 4. ImageNet 256 MF-RAE 1-step samples on classes 13, 14, 94, and 134: snowbird, indigo bird, hummingbird, and crane bird.



Figure 5. ImageNet 256 MF-RAE 1-step samples for various dogs.



Figure 6. ImageNet 256 MF-RAE 1-step samples for class 933: cheeseburger.



Figure 7. ImageNet 256 MF-RAE 1-step samples for class 959: carbonara.



Figure 8. ImageNet 256 MF-RAE 1-step samples for class 947: mushroom.