

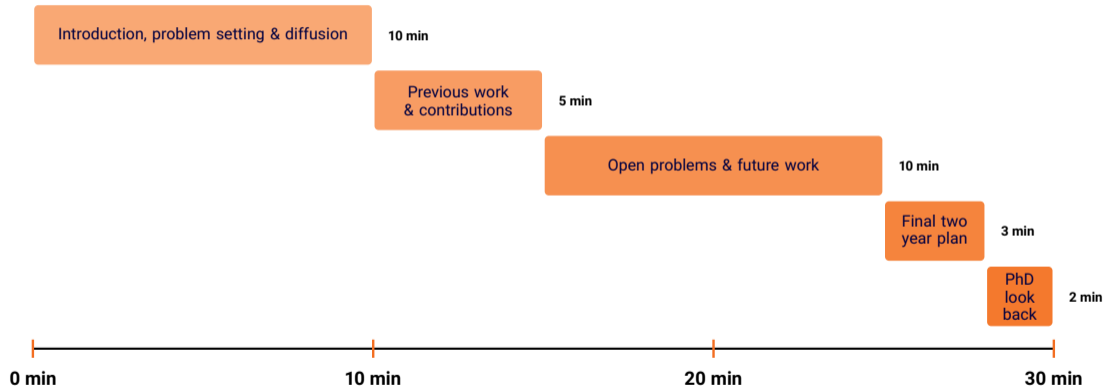
Training-Free Guidance in Diffusion Models: From Classifier Gradients to Future Control Mechanisms

Philipp Vaeth (philipp.vaeth@thws.de)

17-18 June 2026

DMML Workshop Presentation / PhD Exam

Overview and road map



Bayesian setting: sampling from a complex posterior

With $\mathbf{x} \sim p_{\text{data}}$ and \mathbf{y} as a conditioning variable (label, class, measurement,...), we want to draw samples by:

$$\underbrace{p(\mathbf{x} | \mathbf{y})}_{\text{desired posterior}} \propto \underbrace{p_{\theta}(\mathbf{x})}_{\text{learned prior}} \underbrace{p(\mathbf{y} | \mathbf{x})}_{\text{test-time likelihood}} .$$

$$\nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y}) = \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) .$$

Why not learn the posterior directly?

- **Complex conditioning:** The condition \mathbf{y} may be unknown during training or change after training; the likelihood $p(\mathbf{y} | \mathbf{x})$ can be expensive to compute (for example for forward models or simulations).
- **Multiple conditions:** Multiple independent likelihood functions can be applied during sampling:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x} | y_1, \dots, y_k) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \sum_{i=1}^k \nabla_{\mathbf{x}} \log p(y_i | \mathbf{x}) \quad \iff y_1 \perp y_2 \dots \perp y_k | \mathbf{x}$$

- **Flexibility & Efficiency:** The prior $p_{\theta}(\mathbf{x})$ is amortized over many conditions; the likelihood $p(\mathbf{y} | \mathbf{x})$ is amortized over different prior models; plug-and-play.

Diffusion models: score-based overview

Forward SDE

$$d\mathbf{x}_t = f_t(\mathbf{x}_t) dt + g_t d\mathbf{w}_t, \quad \mathbf{x}_0 \sim p_0, \quad \mathbf{x}_T \approx \mathcal{N}(0, \mathbf{I}) .$$

Reverse SDE

$$d\mathbf{x}_t = [f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)] dt + g_t d\bar{\mathbf{w}}_t, \quad T \rightarrow 0, \quad s_\theta(\mathbf{x}, t) \approx \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) .$$

Denosing score matching objective

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}_t} \left[\lambda(t) \|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}_0)\|_2^2 \right], \quad \mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{x}_0, \sigma_t^2 I) .$$

Bayesian guidance: prior score to posterior score

Prior reverse SDE

$$d\mathbf{x}_t = [f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)] dt + g_t d\bar{\mathbf{w}}_t .$$

Score decomposition

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y}) = \underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}_{\text{diffusion prior score}} + \underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} | \mathbf{x}_t)}_{\text{test-time likelihood score}} .$$

Guided posterior reverse SDE

$$d\mathbf{x}_t = [f_t(\mathbf{x}_t) - g_t^2 (s_\theta(\mathbf{x}, t) + \nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} | \mathbf{x}_t))] dt + g_t d\bar{\mathbf{w}}_t .$$

Test-time likelihood: from theory to practical applications

Noise-aware likelihood

$$\nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} \mid \mathbf{x}_t, t) : \phi^* = \arg \max_{\phi} \mathbb{E}_{(\mathbf{x}_0, \mathbf{y}) \sim p_{\text{data}}, \mathbf{x}_t \sim p(\mathbf{x}_t \mid \mathbf{x}_0)} [\log p_\phi(\mathbf{y} \mid \mathbf{x}_t, t)] .$$

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} \mid \mathbf{x}_t) \approx \nabla_{\mathbf{x}_t} \log p(\mathbf{y} \mid D_\theta(\mathbf{x}_t, t) + c),$$

$$\Rightarrow D_\theta(\mathbf{x}_t, t) \approx \mathbb{E}_{p(\mathbf{x}_0 \mid \mathbf{x}_t)} [\mathbf{x}_0] = \frac{\mathbf{x}_t + \sigma_t^2 s_\theta(\mathbf{x}, t)}{\alpha_t} .$$

Plug-and-play

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} \mid \mathbf{x}_t) : \log p_t(\mathbf{y} \mid \mathbf{x}_t) \in \left\{ \log p_\phi(c \mid \cdot), -\ell(\cdot), r(\cdot), -\frac{1}{2\sigma_y^2} \|\mathcal{A}(\cdot) - \mathbf{y}\|_2^2 \right\} .$$

\Rightarrow likelihood model assumptions and setup impact gradient quality.

\Rightarrow requires extensive hyper-parameter searches.

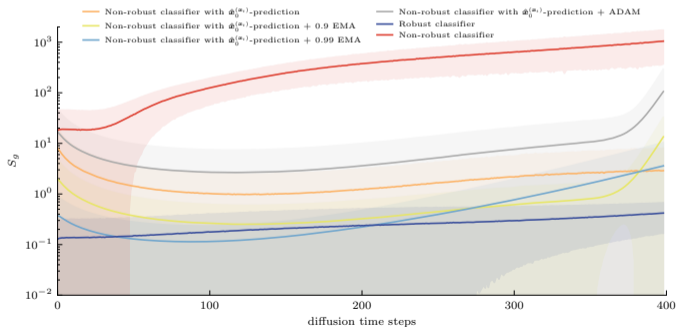
PhD Contributions

	Paper	Research question	TL;DR
I	DBVCE-eval (ECML Workshop 2023)	How do we evaluate (generated) counterfactual explanations?	Compares evaluation criteria for counterfactual explanations.
II	GradCheck (ArXiv 2024)	How to time dynamics of classifier gradients relate to downstream metrics?	Diagnoses when clean classifiers provide useful or unstable forces.
III	nrCG (ECML 2025)	Can stabilization make non-robust classifiers usable?	Extends classifier guidance with denoising and gradient stabilization.
IV	XAIDIFF (ECML Workshop 2025)	How do counterfactual explanations and generative modeling connect?	Frames explanations probabilistically as guided generative sampling.

Case study II/III: non-robust classifier guidance

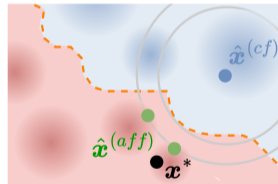
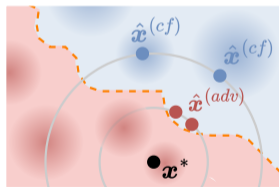
Guidance stability metric

$$S_g(\mathbf{x}_t, \mathbf{x}_{t-1}) = \frac{\|\nabla_{\mathbf{x}_t} f(\mathbf{x}_t) - \nabla_{\mathbf{x}_{t-1}} f(\mathbf{x}_{t-1})\|_2}{\|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2}$$



Case study IV: generative explanations by guided sampling

$$d\mathbf{x}_t = \left[f_t(\mathbf{x}_t) - g_t^2 \left(\underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}_{\text{fidelity}} + \gamma_t \underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} | \mathbf{x}_t)}_{\text{validity}} + \beta_t \underbrace{\left\| \hat{\mathbf{x}}_0^{(\mathbf{x}_t)} - \mathbf{x}^* \right\|_1}_{\text{closeness}} \right) \right] dt + \underbrace{g_t d\bar{\mathbf{w}}_t}_{\text{diversity}}$$



-- Classifier boundary ● Class A Fidelity ● Class B Fidelity

Open questions for future research

1. **Approximation gap** between $\nabla_{\mathbf{x}_t} \log p(c | \mathbf{x}_t)$ and $\nabla_{\mathbf{x}_t} \log p(c | \hat{\mathbf{x}}_0(\mathbf{x}_t, t))$
2. **Likelihood model** properties, differentiability and complexity
3. **Dependent conditions** in multi-condition guidance
4. (optional) **Guidance weight/schedule** γ_t
5. (optional) **Posterior calibration** of diversity, likelihood and prior consistency

Open problem: $\hat{\mathbf{x}}_0$ -prediction / Tweedie Estimation / Mean Guidance

Goal

Approximate the unavailable time-aware likelihood gradient $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} | \mathbf{x}_t)$, interpreting gradient reliability as variance reduction:

$$\text{Var} [\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} | \mathbf{x}_t)] < \text{Var} [\nabla_{\mathbf{x}_t} \log p(c | \hat{\mathbf{x}}_0(\mathbf{x}_t, t) + \Delta_t)] < \text{Var} [\nabla_{\mathbf{x}_t} \log p(c | \hat{\mathbf{x}}_0(\mathbf{x}_t, t))] .$$

Fundamental challenge

Off-manifold clean estimates amplify likelihood-gradient noise:

$$p_0(\cdot) \downarrow \implies \text{Var} [\nabla \log p(\mathbf{y} | \cdot)] \uparrow .$$

RQ1

Can we close the approximation gap $\nabla_{\mathbf{x}_t} \log p(c | \hat{\mathbf{x}}_0(\mathbf{x}_t, t)) \approx \nabla_{\mathbf{x}_t} \log p(c | \mathbf{x}_t)$?

Open problem: likelihood model

Goal

Semantically meaningful guidance gradients from arbitrary likelihood models $\nabla_{\mathbf{x}} \log p_{\phi}(\mathbf{y} | \mathbf{x})$.

Fundamental challenge

Semantic signal in guidance gradient depends on:

- Complexity of the model
- Training objective
- Inductive biases
- Differentiability, discreteness, Lipschitzness

RQ2

Which likelihood model properties lead to stable, semantic and calibrated gradients?

Open problem: multiple not conditionally-independent conditions

Goal

Guide by multiple possibly correlated conditions

$$\nabla_{\mathbf{x}} \log p(\mathbf{x} \mid \mathbf{y}_1, \mathbf{y}_2) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y}_1 \mid \mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y}_2 \mid \mathbf{x}) + \underbrace{\nabla_{\mathbf{x}} \log \frac{p(\mathbf{y}_1 \mid \mathbf{x}, \mathbf{y}_2)}{p(\mathbf{y}_1 \mid \mathbf{x})}}_{\text{dependence residual}} .$$

Fundamental challenge

- What are the implications of dropping the dependence residual?
- How could we estimate or learn the dependence residual?

RQ3

How can we guide by multiple correlated conditions?

Open problem: principled guidance weights

Goal

A principled and efficiently optimized guidance weight:

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t \mid \mathbf{y}) = s_\theta(\mathbf{x}, t) + \gamma(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} \mid \mathbf{x}_t).$$

Fundamental challenge

Optimal hyper-parameters require extensive grid/beam searches but settings vary across:

- downstream tasks and data sets
- likelihood model, prior model and solver
- conditioning targets

RQ4

Can $\gamma(t)$ be derived independent of the a) conditioning variable, b) solver, c) prior and likelihood model, d) data set, and e) tasks?

Open problem: calibrated posterior diversity

Goal

Prior realism, high likelihood and posterior diversity.

$$d\mathbf{x}_t = \left[f_t(\mathbf{x}_t) - g_t^2 \left(\underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}_{\text{prior realism}} + \gamma_t \underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{y} | \mathbf{x}_t)}_{\text{high likelihood}} \right) \right] dt + \underbrace{g_t d\bar{\mathbf{w}}_t}_{\text{diversity}}.$$

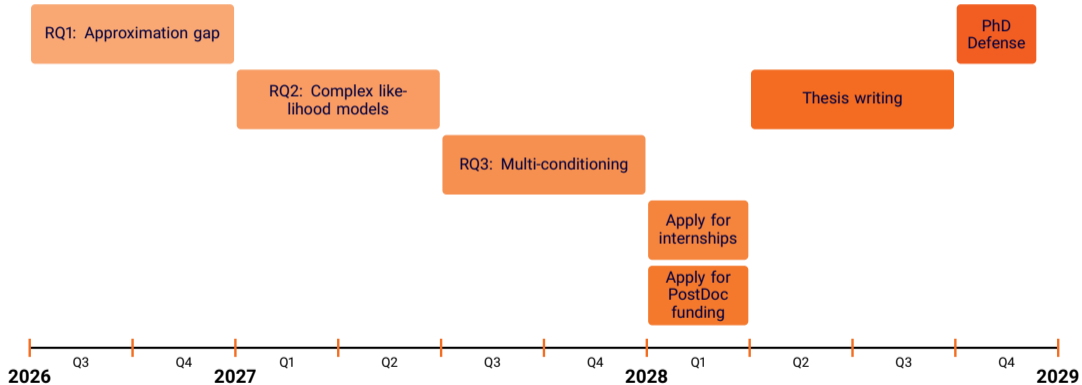
Fundamental challenge

- Prior realism and high likelihood can increase jointly when target labels are well-represented, while often trading-off posterior diversity.
- Guidance has no explicit diversity objective; diversity mainly enters through $g_t d\bar{\mathbf{w}}_t$.

RQ5

Can training-free samplers achieve prior realism and high likelihood while preserving posterior diversity?

Future Plan



PhD retrospective: where I stand now

What went well

1. **Personal growth** from zero math background to generative modeling, Diffusion, SDEs, optimization, (differential) geometry, bounds and proofs.
2. **Research seniority** in the research center, reviewing for conferences, (co-)supervising students, being involved in projects, getting to a broader overview of ML, developing better intuitions.

What I want to improve

1. **Research focus** to go deep when things get hard not pivot (cg > cfg > metrics > memorization > privacy/unlearning > quantization).
2. **Faster prototyping** to quickly judge ideas and either move on or go deep (volumes).
3. **Trust in own ideas** instead of reading 20 papers on the topic and dismissing it because "its solved" (cg).
4. **Pragmatism over perfectionism**

Training-Free Guidance in Diffusion Models: From Classifier Gradients to Future Control Mechanisms

Philipp Vaeth (philipp.vaeth@thws.de)

17-18 June 2026

DMML Workshop Presentation / PhD Exam

References i

- Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *ICLR*, 2023.
- Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. *CVPR Workshop*, 2023.
- Chubin Chen, Jiashu Zhu, Xiaokun Feng, Nisha Huang, Chen Zhu, Meiqi Wu, Fangyuan Mao, Jiahong Wu, Xiangxiang Chu, and Xiu Li. Stochastic self-guidance for training-free enhancement of diffusion models. *ICLR*, 2026.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. *ICLR*, 2023.
- Pierre Del Moral. *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*. Springer, 2004.
- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020.

References ii

- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *NeurIPS*, 2022.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *ICLR*, 2023.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *ICLR*, 2023.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *ICLR*, 2021a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *ICLR*, 2021b.
- Philipp Vaeth, Alexander M. Fruehwald, Benjamin Paassen, and Magda Gregorova. Diffusion-based visual counterfactual explanations – towards systematic quantitative evaluation. *ECML PKDD XKDD Workshop*, 2023. URL <https://arxiv.org/abs/2308.06100>.

References iii

- Philipp Vaeth, Alexander M. Fruehwald, Benjamin Paassen, and Magda Gregorova. Gradcheck: Analyzing classifier guidance gradients for conditional diffusion sampling. *arXiv:2406.17399*, 2024. URL <https://arxiv.org/abs/2406.17399>.
- Philipp Vaeth, Alexander M. Fruehwald, Benjamin Paassen, and Magda Gregorova. Generative example-based explanations: Bridging the gap between generative modeling and explainability. *ECML PKDD XKDD Workshop*, 2025a. URL <https://arxiv.org/abs/2410.20890>.
- Philipp Vaeth, Dibyanshu Kumar, Benjamin Paassen, and Magda Gregorova. Diffusion classifier guidance for non-robust classifiers. *ECML PKDD*, 2025b. URL <https://arxiv.org/abs/2507.00687>.
- Haoxuan Wang, Jinlong Peng, Qingdong He, Hao Yang, Ying Jin, Jiafu Wu, Xiaobin Hu, Yanjie Pan, Zhenye Gan, Mingmin Chi, Bo Peng, and Yabiao Wang. Unicombine: Unified multi-conditional combination with diffusion transformer. *ICCV*, 2025.
- Haotian Ye, Haowei Lin, Jiaqi Han, Minkai Xu, Sheng Liu, Yitao Liang, Jianzhu Ma, James Zou, and Stefano Ermon. Tfg: Unified training-free guidance for diffusion models. *NeurIPS*, 2024.