

Analyzing the Failure Modes of a Hypernetwork-SDE-Flow Architecture in Continual Learning

Vandoanh Van
phamvandoanh9@gmail.com

November 15, 2025

Abstract

Continual learning remains a significant challenge in artificial intelligence, with catastrophic forgetting being a primary obstacle to developing lifelong learning agents. This paper investigates a novel and theoretically promising architecture, the Hypernetwork-SDE-Flow, which combines the task-specific weight generation of hypernetworks, the probabilistic dynamics of Neural Stochastic Differential Equations (SDEs), and the complex distribution modeling of normalizing flows. The primary hypothesis was that this combination would effectively mitigate catastrophic forgetting by isolating task-specific knowledge. However, through a series of rigorous experiments and systematic debugging on a synthetic continual learning benchmark, we demonstrate that this architecture, contrary to our expectations, fails to outperform a simple MLP baseline and suffers from significant catastrophic forgetting. This work provides a detailed analysis of the failure modes, investigating the impact of individual loss components and hyperparameter sensitivity. We conclude by discussing the potential theoretical pitfalls of this combined approach, such as optimization complexity and the challenge of protecting shared parameters, offering valuable insights for future research in the field. Our findings underscore the critical gap between theoretical promise and practical performance in complex continual learning models.

1 Introduction

The ability to learn sequentially from a continuous stream of data is a hallmark of human intelligence but remains a formidable challenge for artificial neural networks. The dominant paradigm of training models on large, static datasets leads to a phenomenon known as catastrophic forgetting [1], where a network's performance on previously learned tasks degrades severely upon learning a new one. Overcoming this limitation is crucial for developing robust, adaptable AI systems capable of lifelong learning.

A variety of methods have been proposed to address catastrophic forgetting, broadly categorized into regularization-based, replay-based, and parameter-isolation methods [2]. Parameter-isolation methods, which aim to assign distinct model parameters to different tasks, are particularly appealing as they offer a clear mechanism for preventing knowledge overwrite. Within this category, hypernetworks [3] have emerged as a powerful approach, utilizing a small network to generate the weights for a larger main network on a per-task basis.

Concurrently, deep generative models have shown promise in learning robust data representations. Neural Stochastic Differential Equations (SDEs) [4] model continuous-time dynamics, offering a powerful framework for capturing the underlying generative processes of data. Normalizing flows [5] provide a mechanism for constructing complex, high-dimensional probability distributions with exact likelihood computation.

In this work, we propose and investigate a novel architecture, the **Hypernet-SDE-Flow**, which synergistically combines these three advanced concepts. Our initial hypothesis was that a hypernetwork could generate task-specific weights for an SDE-based autoencoder, while a normalizing flow learns a structured latent manifold, thereby isolating task knowledge and preventing forgetting. This theoretically grounded approach appeared to be a promising direction for creating a powerful continual learning model.

However, the primary contribution of this paper is not a story of success, but a cautionary tale and a detailed empirical analysis of failure. Despite extensive debugging and hyperparameter tuning, our experiments show that the Hypernet-SDE-Flow model fails to retain knowledge of past tasks, performing significantly

worse than a simple MLP autoencoder baseline. We present a systematic investigation into the model’s behavior, including ablation studies on its complex loss function and an analysis of its performance under different training configurations. Our results highlight the immense challenge of optimizing such complex, multi-component models and suggest that the interaction between shared and task-specific parameters may be more subtle than initially hypothesized. By documenting and analyzing these negative results, we aim to provide valuable lessons for the community and guide future research toward more robust and practical solutions for continual learning.

2 Methodology

The core of our investigation is the Hypernet-SDE-Flow model, a deep generative autoencoder designed for continual learning. The architecture is built upon the principle of parameter isolation, where a hypernetwork generates task-specific parameters for the main encoder and decoder networks. These networks, in turn, learn to map data to and from a structured latent space, whose dynamics are modeled by a Neural SDE and whose distribution is shaped by a normalizing flow.

2.1 Overall Architecture

The model consists of four main components:

1. **A Task-Conditional Encoder:** Maps an input data point x to an initial latent representation z_0 , conditioned on a task ID. This encoder’s weights are dynamically generated.
2. **A Manifold Normalizing Flow:** Transforms the initial latent z_0 into a sample u from a simple base distribution (e.g., a standard Gaussian), allowing for exact likelihood computation.
3. **A Neural Stochastic Differential Equation (SDE):** Models the continuous-time evolution of the latent representation from z_0 to a final state z_T . This captures the generative dynamics of the data.
4. **A Task-Conditional Decoder:** Reconstructs the original data point x from the evolved latent state z_T , also conditioned on the task ID.

2.2 Task-Specific Weight Generation via Hypernetworks

To prevent parameter overwrite between tasks, we employ a ‘HyperNetwork’ that generates the weights and biases for the core layers of the encoder and decoder. For each task t , we associate a learnable task embedding vector e_t . This embedding is fed into the hypernetwork, which consists of small multi-layer perceptrons (MLPs), to output the weight matrices W_t and bias vectors b_t for the ‘DynamicLinear’ layers in the main network. This process ensures that for each task, a functionally distinct encoder and decoder are instantiated, while the hypernetwork itself and other components remain shared across tasks.

2.3 Probabilistic Latent Dynamics with Neural SDEs

We model the generative process in the latent space as a continuous-time stochastic process. The evolution of the latent variable Z is described by the Itô SDE:

$$dZ = f(Z, t)dt + g(Z, t)dW_t \quad (1)$$

where f is the drift function, g is the diffusion function, and W_t is a standard Wiener process. Both f and g are parameterized by neural networks. By solving this SDE (e.g., using an Euler-Maruyama solver), we can transform the initial latent code z_0 into a more expressive representation z_T that captures the underlying dynamics of the data distribution. A separate Neural SDE is used in the decoder to reverse this process.

2.4 Manifold Learning with Normalizing Flows

To enforce a structured geometry on the latent space, we apply a ‘ManifoldNormalizingFlow’ to the initial latent code z_0 . A normalizing flow is a sequence of invertible transformations with easily computable Jacobians. In our model, we use a series of affine coupling layers. This flow transforms the potentially complex distribution of z_0 into a simple, tractable base distribution, $p(u)$, typically an isotropic Gaussian. This allows us to compute the exact log-likelihood of a given data point, which forms a crucial component of our loss function:

$$\log p(x) = \log p(u) + \log |\det(J)| \quad (2)$$

where J is the Jacobian of the entire transformation from x to u . This encourages the model to learn a well-formed latent manifold, which is hypothesized to aid in task separation.

2.5 Loss Function

The model is trained end-to-end by minimizing a composite loss function designed to balance data reconstruction with continual learning objectives:

$$L_{\text{total}} = L_{\text{recon}} + \beta \cdot L_{\text{nll}} + \gamma \cdot L_{\text{replay}} + \delta \cdot L_{\text{sep}} \quad (3)$$

- L_{recon} : A standard mean squared error (MSE) loss to ensure the autoencoder can accurately reconstruct the input data.
- L_{nll} : The negative log-likelihood derived from the normalizing flow, which encourages the model to learn a good probabilistic representation. The weight β is annealed during training.
- L_{replay} : A replay-based loss calculated on a small episodic memory of samples from past tasks, forcing the model to retain knowledge of previous data distributions.
- L_{sep} : A task separation loss that encourages the task embedding vectors e_t for different tasks to be distant from each other in the embedding space.

3 Experiments and Analysis

We conduct a series of experiments to empirically evaluate the performance of the Hypernet-SDE-Flow model. Our primary goal is to assess its ability to mitigate catastrophic forgetting in a sequential learning setting.

3.1 Experimental Setup

Benchmark Dataset: We use a synthetic benchmark consisting of 5 distinct tasks. Each task is a dataset of 300 samples, where each sample is a 3-dimensional vector generated by a different non-linear function. This controlled environment allows for a clear and interpretable evaluation of forgetting.

Baseline Model: We compare our model against a standard **MLP Autoencoder**. This baseline has a similar number of parameters in its main encoder/decoder structure but lacks any specific mechanism for continual learning. It is trained sequentially on the tasks, with its weights being overwritten at each step.

Training Details: Both models are trained for 300 epochs per task using the AdamW optimizer with a learning rate of 5e-4. The Hypernet-SDE-Flow model was trained using the full composite loss function described in Section 2.5 after confirming that simpler versions of the loss did not resolve the core issues.

3.2 Primary Results: Catastrophic Forgetting Analysis

After training on all 5 tasks sequentially, we evaluate the performance of both models on all tasks. The reconstruction Mean Squared Error (MSE) is used as the evaluation metric. Table 1 and Table 4 show the final results.

Table 1: Final Reconstruction MSE Matrix. Shows the MSE on task T_j after the model has finished training on task T_i .

Table 2: Baseline MLP					
Trained On	Test T0	Test T1	Test T2	Test T3	Test T4
Task 0	0.0044	-	-	-	-
Task 1	0.0053	0.0001	-	-	-
Task 2	0.0093	0.0122	0.0000	-	-
Task 3	0.0069	0.0114	0.0008	0.0000	-
Task 4	0.0060	0.0046	0.0011	0.0005	0.0000

Table 3: Hypernet-SDE-Flow					
Trained On	Test T0	Test T1	Test T2	Test T3	Test T4
Task 0	0.0083	-	-	-	-
Task 1	0.2166	0.0090	-	-	-
Task 2	0.1267	0.5295	0.0035	-	-
Task 3	1.4147	0.3772	0.9962	0.0016	-
Task 4	0.3905	0.5956	0.6441	0.0926	0.0087

Table 4: Forgetting Metric on Task 0. We define the forgetting metric as the increase in loss on Task 0 after training on all subsequent tasks.

Model	Loss on T0 (after T0)	Loss on T0 (after T4)	Forgetting (Increase)
Baseline MLP	0.0044	0.0060	+0.0016
Hypernet-SDE-Flow	0.0083	0.3905	+0.3822

The results clearly indicate that the Hypernet-SDE-Flow model suffers from severe catastrophic forgetting. For instance, the MSE on Task 0 balloons from 0.0083 to 0.3905 after training on subsequent tasks. In contrast, the baseline MLP’s performance on Task 0 remains remarkably stable.

The Hypernet-SDE-Flow model forgets approximately 238 times more than the simple MLP baseline, a result that directly contradicts the initial hypothesis.

3.3 Debugging Analysis

To understand this failure, we performed several diagnostic experiments:

1. **Loss Simplification:** We ran an experiment where the L_{replay} and L_{sep} were removed to isolate the core architecture. The catastrophic forgetting problem persisted, indicating the issue is not with the auxiliary losses but with the fundamental model structure.
2. **Increased Forgetting Penalties:** We ran another experiment where the weights of the L_{replay} and L_{sep} were significantly increased. This also failed to mitigate forgetting, suggesting the optimization process is unable to balance the competing objectives effectively.

4 Discussion

The stark contrast between the theoretical elegance of the Hypernet-SDE-Flow and its empirical failure prompts a deeper discussion into why it did not work.

The Illusion of Parameter Isolation: The core assumption was that generating task-specific weights via a hypernetwork would isolate knowledge. However, our results suggest this isolation was insufficient. A

critical oversight is the substantial number of **shared parameters** that remain. The hypernetwork itself, the task embedding vectors, the entire normalizing flow, and the neural networks within the SDEs are all shared across tasks. When training on a new task, the gradients flow through these shared components, inevitably altering them in a way that is detrimental to past tasks. The model, in effect, learns to be a good "generalist" for the most recently seen data, overwriting the specialized functions of these shared components.

Optimization Complexity: By combining four distinct and complex components (Autoencoder, Hypernet, SDE, Flow), we create an extraordinarily complex and non-convex optimization landscape. The composite loss function, with its multiple competing objectives and annealed hyperparameters, is notoriously difficult to optimize. It is highly probable that the optimizer converges to a sharp local minimum that is good for the current task but poor for previously seen tasks, rather than finding a flatter, more generalizable solution.

Lessons Learned: This investigation provides several key takeaways for the continual learning community. Firstly, architectural complexity is not a panacea. While combining powerful theoretical concepts is appealing, it can lead to models that are practically untrainable or behave in unexpected ways. Secondly, true parameter isolation must be more explicit and complete than what was achieved here. Future work could explore methods that enforce stricter parameter disjointedness, perhaps through architectural constraints or targeted regularization.

5 Conclusion

In this paper, we presented a systematic investigation of a novel continual learning architecture, the Hypernet-SDE-Flow. Despite its strong theoretical grounding, our comprehensive experiments and debugging efforts led to an unambiguous conclusion: the model fails to mitigate catastrophic forgetting and is significantly outperformed by a simple baseline. We have provided a detailed analysis of this failure, hypothesizing that the root causes lie in incomplete parameter isolation and insurmountable optimization challenges. While a negative result, our work serves as a valuable case study, highlighting the pitfalls of architectural complexity and providing a clear empirical data point that can help guide the community towards more robust and effective solutions for the grand challenge of lifelong learning.

References

- [1] McCloskey, M., & Cohen, N. J. (1989). Catastrophic interference in connectionist networks: The sequential learning problem. In *The psychology of learning and motivation*, Vol. 24, pp. 109-165. Academic Press.
- [2] Parisi, G. I., Kemker, R., Part, J. L., Kanan, C., & Wermter, S. (2019). Continual lifelong learning with neural networks: A review. *Neural Networks*, 113, 54-71.
- [3] Ha, D., Dai, A., & Le, Q. V. (2016). Hypernetworks. In *International Conference on Learning Representations (ICLR)*.
- [4] Li, X., Wong, T. K., Chen, R. T., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [5] Rezende, D. J., & Mohamed, S. (2015). Variational inference with normalizing flows. In *International Conference on Machine Learning (ICML)*.