

Latent-Autoregressive GP-VAE Language Model

Yves Ruffenach

Conservatoire National des Arts et Métiers

yves.ruffenach.auditeur@lecnam.net

yves@ruffenach.net

ORCID: [0009-0009-4737-0555](https://orcid.org/0009-0009-4737-0555)

Abstract

We investigate a fully Latent AutoRegressive scheme based on a Gaussian Process (GP) integrated into a Variational Autoencoder (VAE). In this setting, sequential dynamics are transferred from the observation space to a continuous latent space, while linguistic generation remains parallel through a non-autoregressive decoder. We present a complete methodological formulation, including a causal GP prior, a structured amortized posterior, and a training protocol based on a regularized ELBO. Empirical evaluation, conducted within a deliberately constrained proof-of-concept (POC) framework, shows that the model can be trained stably and that the sequential and parallel sampling variants exhibit consistent behavior. Overall, the results suggest that part of the temporal structure in a language model can be supported by the probabilistic geometry of the latent space rather than by explicit neural operations.

Keywords. Gaussian Process VAE; Sequential models; Latent autoregression; Language modeling; Reasoning; Bayesian generative models; Deep learning; Variational autoencoders; Gaussian processes.

1 Introduction

Recent advances in language models are no longer driven solely by increasing their size or by refining their internal mechanisms. Mixtures of Experts, large-scale architectures, and extensive fine-tuning have demonstrated their effectiveness, yet the collective experience of the field now highlights a central observation: the strongest models are those that reason. They combine multiple perspectives, explore alternative inference chains, interact with one another, and rely on distributed update mechanisms within a broader ecosystem of models.

This reflects a much more pragmatic intuition: coherent language generation requires an underlying coherent internal dynamic. A model that is able to organize a stable and structured latent evolution is, in practice, a model capable of producing stable and structured linguistic outputs. From this perspective, the link between reasoning and language does not arise from symbolic manipulation alone, but from the ability to maintain a consistent internal trajectory that guides generation.

This work adopts that perspective. Rather than taking language as the model’s starting point, we begin with a reasoning dynamic: a continuous internal structure that precedes, guides, and constrains symbolic generation. This idea builds on a series of works showing that Variational Autoencoders (VAEs) provide a latent expressiveness particularly well suited to one-shot learning, rapid adaptation, and the geometric shaping of continuous representations—properties that are difficult to obtain in strictly autoregressive architectures.

We therefore explore the following hypothesis: the sequential dynamics of a language model can be shifted from the observable space into the latent space. Instead of enforcing causality on tokens through explicit autoregression or large-scale self-attention, we assign temporal dependence to an

analytical covariance defined by a Gaussian Process (GP). Causality thus becomes a mathematical property of the latent space rather than a neural loop applied to the observations.

It is essential to emphasize that the autoregressivity considered in this work operates exclusively within the latent space. It should not be confused with the classical token-level autoregression used in models such as Transformers, which directly conditions on observed sequences. This distinction is fundamental: our model explores an internal latent causality rather than an explicit sequential dependence between tokens.

This perspective leads to the model studied in this paper: a GP-VAE equipped with a purely latent autoregressive mechanism, in which

- the GP prior enforces internal continuity and causality,
- the encoder learns a structured amortized posterior,
- the decoder remains fully parallel.

In this framework, language is no longer generated through an autoregressive roll-out but as the projection of a continuous latent trajectory. The following chapters formalize this paradigm, present the methodology enabling it to scale, and assess its empirical validity within a controlled proof-of-concept setting.

Finally, note that a correlated Gaussian Process imposes a metric geometry but no inherent directionality: two points equally close under the kernel are treated as equivalent. Such symmetry can introduce ambiguities in a sequential task like language modeling. The introduction of purely latent autoregressivity resolves this ambiguity: by forcing each z_t to depend on the history $z_{<t}$, it provides the GP with an explicit directional structure and allows the model to disambiguate latent contributions that would otherwise be interchangeable. The resulting oriented latent dynamics combine the GP’s correlated geometry with a coherent causal progression.

2 Related Work

2.1 Variational and Deep Generative Models

Since the introduction of the Variational Autoencoder (VAE) by Kingma & Welling [10], variational generative models have become a central framework for latent probabilistic modeling. A VAE links a continuous latent space z to a complex observation space x through an encoder $q_\phi(z | x)$ and a decoder $p_\theta(x | z)$, optimized via the Evidence Lower BOund (ELBO):

$$\mathcal{L} = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x | z)] - D_{\text{KL}}(q_\phi(z | x) \parallel p(z)). \quad (1)$$

The KL term controls the regularization of the latent space. Among notable variants, the β -VAE [6] rescales D_{KL} to encourage disentanglement of latent factors — meaning the separation of hidden dimensions that structure internal representations. Conditional VAEs extend learning to external context; hierarchical models [19] and sequential VAEs [3, 9] aim to represent multi-scale temporal dependencies.

All such approaches seek to capture transferable latent regularities beyond the observed data, which is beneficial for few-shot and one-shot learning [18].

2.2 Gaussian Processes and Bayesian Deep Learning

A Gaussian Process (GP) [17] defines a distribution over functions:

$$f(t) \sim \mathcal{GP}(m(t), k(t, t')), \quad (2)$$

where m is the mean and k the covariance kernel.

Each observation refines the matrix K , making the GP fundamentally nonparametric: its complexity grows with the number of data points. The combination of GPs with neural networks via *Deep Kernel Learning* [22] has enabled models that blend neural flexibility with Bayesian structure.

Used as a latent prior, the GP yields GP-VAEs, where the covariance correlates latent variables z_t and $z_{t'}$. Casale et al. [1] introduced the *Gaussian Process Prior VAE*, structuring the latent space; Fortuin et al. [2] applied it to temporal imputation; Pearce et al. [15, 14] developed additional GP-based Bayesian autoencoders.

Such models learn not independent latent points but coherent latent trajectories, whose continuity arises from the kernel $k_\psi(t, t')$. However, the induced correlations remain symmetric and do not encode explicit temporal causality.

2.3 Toward Latent Autoregression and One-Shot Generalization

Rezende et al. [18] showed that deep generative models can produce coherent samples from a single example (*one-shot generalization*). These results motivated approaches where a correlated prior enhances contextual coherence under limited data.

More recent work, such as Wang et al. [21], has shown that dense Gaussian Processes embedded in deep networks substantially improve few-shot performance.

The central insight is that a GP, by inducing latent correlations, can also support a latent sequential dynamic. Causal dependence can thus be transferred from the observable space to the continuous latent space.

The conditional factorization of a causal GP is:

$$p(z_t \mid z_{<t}) = \mathcal{N}\left(k_{12}^\top K_{11}^{-1} z_{<t}, K_{22} - k_{12}^\top K_{11}^{-1} k_{12}\right), \quad (3)$$

which naturally allows parallel decoding:

$$p_\theta(x \mid z_{1:L}). \quad (4)$$

Thus, sequential memory is provided by GP covariance rather than an explicit neural recurrence.

2.4 Latent Sequential Dynamics and the Novelty of Purely Latent Autoregression

Sequential latent-variable models (dynamic VAEs, latent state-space models, RNN-VAEs, etc.) explore various temporal dependencies, but none explicitly formalizes what we call *purely latent autoregression*:

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t \mid z_{<t}), \quad (5)$$

in which each z_t conditions on the full latent history $z_{<t}$.

Earlier models—especially Fraccaro et al. [3] and Karl et al. [9]—use Markovian transitions $p(z_t \mid z_{t-1})$. Kingma et al. [11] introduced Inverse Autoregressive Flows, which operate over latent densities but do not represent a causal temporal dynamic.

More recent studies, such as Zhou et al. [24] and Klushyn et al. [12], employ expressive continuous latent-state models but still do not formalize full causal dependence.

To the best of our knowledge, no prior work combines simultaneously:

- a correlated GP prior,
- causal latent sequential generation $p(z_t \mid z_{<t})$,
- a fully non-autoregressive decoder.

We propose this combination as an original contribution: a GP-VAE with purely latent autoregression, where causality is embedded directly into the probabilistic structure of the latent process. The formulation is:

$$p(z_{1:L}) = \prod_{t=1}^L \mathcal{N}(z_t; f_\psi(z_{<t}), \Sigma_\psi(z_{<t})), \quad (6)$$

which provides analytical temporal continuity and internal Bayesian coherence. This structure encompasses Markov and AR(p) processes as special cases, and offers an analytical alternative to attention-based modeling by shifting the causal mechanism into the latent space.

Novelty clarification. While latent sequential models have been explored in prior work (e.g., Fraccaro et al., Karl et al.), these approaches rely on *parametric* state transitions—typically Markovian forms

$$p(z_t | z_{t-1}),$$

whose temporal structure is learned through neural weights. The contribution of the present work is not an increase in the order of such transitions, but a *change in the nature of the causal mechanism itself*. The full factorization

$$p(z_t | z_{<t})$$

is derived analytically from Gaussian conditioning within a correlated Gaussian Process prior, making temporal dependence a *probabilistic-geometric* property rather than a learned recurrence.

To our knowledge, no prior work combines simultaneously:

- a correlated Gaussian Process prior;
- an explicitly causal latent factorization

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t | z_{<t});$$

- a fully non-autoregressive decoder.

This combination shifts sequential dynamics from the symbolic level to the latent space, constituting a conceptual—rather than incremental—extension of existing latent sequential models.

2.5 Current Limitations and Conceptual Positioning

Existing GP-VAEs face three major limitations:

1. the absence of temporal causality in the standard GP;
2. the cubic complexity $O(L^3)$ of covariance inversion;
3. the lack of explicit sequential dynamics.

The proposed model—a GP-VAE with purely latent autoregression—addresses these limitations: the GP ensures global probabilistic coherence; the factorization $p(z_t | z_{<t})$ introduces explicit causal direction; and the parallel decoder enables simultaneous generation.

This structure replaces neural memory with analytical covariance, reducing computational cost while maintaining semantic coherence. It provides a conceptual alternative to large autoregressive models: linguistic coherence emerges from probabilistic continuity rather than symbolic recurrence.

This constitutes a new methodological contribution not previously implemented in the literature.

2.6 Implementation and Associated Tools

The implementation relies on GPyTorch [4]. The BBMM algorithm (*Blackbox Matrix–Matrix Multiplication*) enables covariance inversion in $O(L^2)$ via vectorized conjugate gradients. This framework supports kernel hyperparameter optimization and scalable GPU training.

The kernels used (RBF, Matérn [5], Spectral Mixture [23]) induce distinct latent geometries, enabling exploration of internal temporal continuity or periodicity.

2.7 Summary

This chapter positions the proposed model at the intersection of three major research lines:

- variational autoencoders [10, 6, 19, 3],
- correlated latent models of the GP-VAE family [1, 2, 15, 7],
- one-shot generalization approaches [18, 21].

Purely latent autoregression extends GP-VAEs by introducing explicit temporal causality alongside a non-autoregressive decoder. It shifts causality from the symbolic level to a continuous probabilistic space, enabling more compact, stable, and interpretable language models.

While previous approaches have shown that Gaussian Processes provide useful statistical continuity, none has yet formalized fully latent, fully sequential causality. The purely latent autoregression introduced here extends this direction by transferring temporal dependence from the observation space to the latent space.

The next chapter presents the methodology that implements this model: reducing computational complexity, structuring the encoder, defining the causal latent prior, designing the parallel decoder, and establishing the experimental protocol. The goal is to demonstrate how this theoretical framework becomes a scalable system capable of handling long sequences while maintaining the Bayesian coherence of the latent process.

3 Proposed Models – Methodology

3.1 Reducing Computational Complexity

One of the main challenges in integrating a Gaussian Process (GP) into a sequential model lies in its computational cost. Evaluating a temporal kernel over T steps yields a covariance matrix $K \in \mathbb{R}^{T \times T}$ whose exact inversion requires $O(T^3)$. This cubic dependence makes naïve GP usage impractical for long sequences unless dedicated approximations are used (inducing points, conjugate gradients, structural factorizations).

Among such methods, the inducing-point approximation [20] reduces inference to a cost of $O(TM^2)$ by selecting $M \ll T$ representative latent points $u = \{u_1, \dots, u_M\}$. GPYTORCH [4], on its side, provides the BBMM (Blackbox Matrix–Matrix Multiplication) algorithm, which approximates the inversion of K using vectorized conjugate gradients, with an effective complexity close to $O(T^2)$. Combined with a spatio-temporal kernel factorization of the form

$$K = K_t \otimes K_s,$$

this approach extends the applicability of GPs to sequences far longer than those manageable through exact inversion.

However, the most decisive gain does not come solely from these numerical optimizations, but from a shift in the level of representation. In a classical autoregressive model (e.g., a Transformer), sequential dependence is carried by the N observable tokens, with each causal-attention layer requiring $O(N^2d)$ operations. For example, a sequence of $N = 1024$ tokens with a hidden

dimension $d = 768$ demands tens of millions of operations per layer, along with substantial intermediate storage for the Q , K , and V matrices.

In the model we propose, this dependence is transferred to a compact latent sequence $z_{1:L}$, typically 8 to 16 times shorter than the token sequence, with latent vectors $z_t \in \mathbb{R}^{d_z}$ where $d_z \ll d$. The temporal dependence in the latent then has cost

$$O(L^2 d_z),$$

representing one to two orders of magnitude less than a model operating directly on observations.

This estimate corresponds to the effective cost in our implementation, which systematically uses GP approximations (BBMM, inducing points) to maintain quadratic inference. Exact covariance inversion, at $O(L^3)$, is not used in practice.

Thus, the reduction is not merely an approximation of the GP, but the result of a change in computational topology: sequential dynamics are encoded in the covariance of the Gaussian Process rather than recomputed via self-attention at every layer.

In other words, the model replaces a dense attention graph with a factorizable covariance structure whose temporal continuity is analytically defined by the kernel. Decoding can then be performed in a single parallel pass:

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t \mid z_{<t}), \quad p(x \mid z_{1:L}) \text{ (parallel),}$$

while a classical autoregressive model requires evaluating

$$p(x_t \mid x_{<t}) \quad \text{for each } t.$$

This shift from the symbolic level to the latent level has two direct consequences: (i) computation becomes fully vectorized rather than sequential; (ii) backpropagation no longer follows a long causal chain, since the temporal dynamics are integrated into the GP covariance.

In summary, the complexity reduction observed in the GP-VAE does not arise solely from numerical optimizations, but from a conceptual relocation of causality: sequential dynamics are no longer carried by an explicit neural memory (RNN, Transformer), but by a correlated probabilistic structure. This change in perspective, combined with GP approximations maintaining quadratic cost, enables Bayesian architectures capable of modeling long sequences while preserving robust temporal and semantic coherence.

3.2 Architecture and Role of the Decoder

The overall architecture of the proposed model relies on a strict functional separation between three components: the amortized encoder, the correlated latent block (GP-AR), and the generative decoder. This separation, typical of variational models, becomes structurally essential here: it isolates the sequential dynamics within the latent space and delegates the reconstruction of observations to a parallel decoder. Unlike classical autoregressive models, where temporal dependence appears directly at the token level, the temporality of the GP-VAE is carried entirely by the latent sequence $z_{1:L}$.

The amortized encoder maps observations into the latent space. It parameterizes an approximate posterior distribution

$$q_\phi(z_{1:L} \mid x_{1:N}), \tag{7}$$

typically Gaussian, whose mean and variance are produced by a hierarchical convolutional or temporal network, such as a Temporal Convolutional Network (TCN), i.e., a causal dilated-convolution architecture suited to sequential dependencies. Its role is purely inferential: to produce an amortized estimate of the latent distribution from input data. Choosing a convolutional

encoder rather than a full Transformer limits training cost, since most of the sequential coherence is enforced by the GP prior.

The correlated latent block forms the core of the model. It defines an autoregressive Gaussian Process over the latent sequence:

$$p_{\theta}(z_{1:L}) = \prod_{t=1}^L p_{\theta}(z_t \mid z_{<t}), \quad p_{\theta}(z_t \mid z_{<t}) = \mathcal{N}(m_t, \Sigma_t), \quad (8)$$

where the moments (m_t, Σ_t) are obtained through Gaussian conditioning on previous steps. The kernel $k_{\psi}(t, t')$ encodes the covariance structure of the process: a squared-exponential or Matérn kernel [5] enforces temporal continuity, while a spectral kernel [23] captures latent periodicities useful for recurrent linguistic structures.

This block plays a role analogous to that of memory mechanisms in sequential neural architectures, but its dynamics are strictly probabilistic: temporal dependence is modeled analytically by the covariance K , not parametrically through learned transition weights.

Finally, the generative decoder maps the latent sequence to concrete observations (tokens, spectrograms, pixels, etc.). It parameterizes the distribution

$$p_{\theta}(x_{1:N} \mid z_{1:L}) = \prod_{n=1}^N p_{\theta}(x_n \mid z_{1:L}), \quad (9)$$

allowing fully parallel generation: all tokens are produced simultaneously from the complete latent trajectory.

By construction, this formulation is compatible with multiple decoder families: convolutional networks, lightweight Transformers, or restricted-attention architectures. The latent factorization guarantees strong modularity: the decoder can be replaced without altering the internal dynamics, since the latter are entirely carried by the GP-AR.

The model thus differs from token-by-token autoregressive architectures, where each new symbol depends on the history $x_{<t}$. Here, temporal dependence has already been integrated into $z_{1:L}$; the decoder merely applies a global latent–observable mapping. Conceptually, this corresponds to the following distinction:

$$(\text{LLM}) \quad p(x) = \prod_t p(x_t \mid x_{<t}) \quad \text{vs.} \quad (\text{GP-VAE}) \quad p(x) = \int p(x \mid z) p(z) \, \mathrm{d}z. \quad (10)$$

Two levels of coherence can thus be distinguished:

- a local coherence, resulting from the proximity between z_t and z_{t-1} , which ensures smooth transitions between consecutive units;
- a global coherence, imposed by the covariance of the Gaussian Process, which governs higher-level regularities (theme, tone, argumentative structure).

Thus, the GP acts as a continuous regulator of global structure, while the decoder imposes the local constraints of the language. In summary, the strength of the model does not lie in decoder complexity, but in the probabilistic structure of the latent space. The decoder becomes an interpreter of an already coherent dynamic: it transforms a continuous trajectory into a symbolic sequence. This decentralization of causality—from tokens to latents—is the methodological pivot of the model, enabling parallel, stable, and Bayesian generation without explicit attention on the observations.

3.3 Purely Latent Autoregression

The notion of purely latent autoregression is one of the key contributions of the proposed model. It refers to a sequential dynamic entirely confined to the latent space—that is, a causal temporal dependence that unfolds between the hidden variables $z_{1:L}$, independently of any observation x . This approach differs both from simple Markov processes (limited to local memory) and from classical neural autoregressive models (where causality is placed directly on observed tokens).

3.3.1 General Principle

A purely latent autoregressive model constructs the sequence of internal states as a full probabilistic chain:

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t \mid z_{<t}), \quad (11)$$

where each conditional distribution is Gaussian:

$$p(z_t \mid z_{<t}) = \mathcal{N}(\mu_t, \Sigma_t), \quad (12)$$

$$\mu_t = k_{(t,<t)}^\top K_{(<t,<t)}^{-1} z_{<t}, \quad (13)$$

$$\Sigma_t = k_{tt} - k_{(t,<t)}^\top K_{(<t,<t)}^{-1} k_{(t,<t)}. \quad (14)$$

These expressions follow from classical Gaussian conditioning and show that the model’s internal dynamics are governed not by a recurrent network but by the covariance structure of the Gaussian Process. Each new latent state is sampled from a predictive distribution updated by previous latents, ensuring smooth and analytical temporal continuity.

This structure can be seen as a Bayesian form of autoregression: memory of the past is transmitted not through neural weights but through probabilistic correlations. Generation occurs entirely in the latent space: the model first constructs a coherent trajectory $z_{1:L}$, and the decoder subsequently maps this trajectory into an observable sequence. Unlike an RNN or a Transformer, no symbolic feedback enters the generation loop: the model never conditions on its own text outputs—hence the qualifier *purely* latent.

3.3.2 Difference from Markov Approaches and Token-Based Methods

It is crucial to distinguish this structure from Markovian or classical autoregressive models. A Markov latent model defines a local dependence:

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t \mid z_{t-1}), \quad (15)$$

that is, a first-order memory. This assumption simplifies inference but limits the temporal range of dependencies: it is well suited for physical trajectories or continuous signals with short-range correlations, but insufficient to capture the syntactic or thematic patterns of language.

Conversely, an autoregressive observed model, for instance GPT or a causal Transformer, operates directly on tokens:

$$p(x) = \prod_{t=1}^N p(x_t \mid x_{<t}), \quad (16)$$

imposing explicit causality but at the cost of a combinatorial context expansion and slow sequential generation. In such models, temporal dependence is symbolic: each linguistic decision conditions on all preceding ones, making the model sensitive to local errors and incompatible with parallel generation.

The latent autoregressive GP-VAE sits in between: it preserves full autoregressive causality (each z_t depends on all $z_{<t}$) while operating in a continuous space where covariance imposes

smooth, probabilistic coherence. The flow of dependence remains causal, but it is decoupled from the symbolic level—syntax and grammar emerge afterward through the decoder.

3.3.3 Intuitive Interpretation and Scope

Operationally, purely latent autoregression appears as a sequential unfolding of Gaussian conditioning: at each step t , the model updates the distribution of z_t based on the latent history $z_{<t}$, according to

$$z_t \sim \mathcal{N}(k_{(t,<t)}^\top K_{(<t,<t)}^{-1} z_{<t}, k_{tt} - k_{(t,<t)}^\top K_{(<t,<t)}^{-1} k_{(t,<t)}). \quad (17)$$

This iteration produces a correlated and causal chain in the latent space, where memory of the past is transmitted through the covariance of the process. The major difference from neural autoregression lies in the analytical nature of the conditioning: temporal coherence emerges from a learned covariance structure rather than from sequential weights.

Conceptually, this amounts to shifting the temporal structure of language into a more abstract internal space. The latent sequence encodes the underlying semantic dynamics—theme transitions, argumentative continuity, syntactic dependencies—before any observable emission. The decoder merely projects this internal dynamic onto the linguistic surface, reducing variance and enhancing stylistic coherence.

3.3.4 Formal Definition

We say that a model exhibits latent autoregression when the joint distribution of hidden variables factorizes causally as

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t \mid z_{<t}), \quad (18)$$

with

$$p(z_t \mid z_{<t}) = f_\theta(z_{<t}) + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \Sigma_t). \quad (19)$$

The function f_θ may be a Gaussian Process, a masked flow, or a small latent Transformer, but in all cases it depends only on past states. When the dependence is restricted to z_{t-1} , one recovers a first-order Markov transition; when it extends to the full history $z_{<t}$, we obtain full latent autoregression.

This formulation ensures complete internal causality, distinct from the external causality of observable models. It preserves the analytical regularity of the GP while incorporating the directionality characteristic of sequential models. The result is a causally oriented latent dynamic, compatible with amortized inference and parallel VAE decoding.

3.3.5 Summary

Purely latent autoregression extends and unifies several paradigms:

- it generalizes local Markovian dependence into global causality;
- it replaces attention mechanisms with analytical covariance;
- it preserves the Bayesian factorization that underlies VAE stability.

The model no longer learns an explicit sequential memory, but a continuous dependence structure within a probabilistic latent space. This reformulation shifts the complexity of language modeling from the token level to latent geometry—a conceptual translation that underpins the remainder of the methodology.

3.4 Probabilistic Formulation and Objective Function

The complete formulation of the model relies on the standard variational framework, enriched with a correlated and causal latent structure. The autoregressive GP-VAE defines a joint distribution over observations x and latent variables z :

$$p_\theta(x, z) = p_\theta(x | z) p_\theta(z), \quad (20)$$

where $p_\theta(z)$ is a correlated, causal latent prior, and $p_\theta(x | z)$ the generative model (the decoder).

Training consists in maximizing the marginal likelihood of the data:

$$\log p_\theta(x) = \log \int p_\theta(x, z) dz, \quad (21)$$

which is not analytically tractable, hence the classical variational approximation:

$$\log p_\theta(x) \geq \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x | z)] - D_{\text{KL}}(q_\phi(z | x) \| p_\theta(z)), \quad (22)$$

the Evidence Lower BOund (ELBO), where D_{KL} denotes the Kullback–Leibler divergence measuring the discrepancy between the approximate posterior and the latent prior.

3.4.1 Correlated Latent Prior

The latent prior plays a central role. Instead of imposing a factorized independence between components of z (as in a standard VAE), we use an autoregressive Gaussian Process that explicitly encodes temporal correlations and causality:

$$p_\theta(z_{1:L}) = \prod_{t=1}^L \mathcal{N}(z_t; k_{(t, < t)}^\top K_{(< t, < t)}^{-1} z_{< t}, k_{tt} - k_{(t, < t)}^\top K_{(< t, < t)}^{-1} k_{(t, < t)}). \quad (23)$$

The covariance matrix K is defined by a differentiable kernel $k_\psi(t, t')$ parameterized by hyperparameters ψ (temporal scale, variance, etc.). This construction guarantees internal continuity and causality: each z_t depends on its predecessors through Gaussian conditioning, and the whole forms a coherent latent trajectory.

3.4.2 Amortized Posterior

The variational posterior $q_\phi(z | x)$ is chosen to be factorized in form while depending implicitly on the context encoded by the encoder:

$$q_\phi(z | x) = \prod_{t=1}^L \mathcal{N}(z_t; \mu_\phi(x_{\leq t}), \sigma_\phi(x_{\leq t})^2 I). \quad (24)$$

This form enables amortized inference: the encoder learns a direct mapping $x \mapsto (\mu, \sigma)$, making posterior evaluation independent of iterative optimization. The use of convolutional or TCN layers promotes parallelization while preserving temporal structure.

3.4.3 Generative Decoder

The decoder $p_\theta(x | z)$ implements the latent–observation projection. In our framework, it is non-autoregressive, meaning that it generates all positions in parallel from the full latent sequence. Each observation is typically modeled as:

$$p_\theta(x | z) = \prod_{n=1}^N \mathcal{N}(x_n; f_\theta(z)_n, \sigma_x^2 I) \quad \text{or} \quad p_\theta(x | z) = \prod_{n=1}^N \text{Cat}(f_\theta(z)_n), \quad (25)$$

depending on whether the data are continuous (series, images) or discrete (tokens).

The network f_θ may be a convolutional architecture or a lightweight Transformer applied to the latents; it maps the correlated regularities into observable structures (grammar, tone, style).

3.4.4 Full Objective Function

Combining these elements, the ELBO of the model becomes:

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x | z)] - D_{\text{KL}}(q_\phi(z | x) \| p_\theta(z)). \quad (26)$$

The first term maximizes reconstruction likelihood (decoder fidelity), while the second regularizes the latent structure by bringing the posterior closer to the GP-AR prior. In practice, the KL decomposes analytically into a sum of weighted Gaussian divergences, and can be modulated by a coefficient β [6] to adjust the trade-off between reconstruction and regularization:

$$\mathcal{L}_\beta = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x | z)] - \beta D_{\text{KL}}(q_\phi(z | x) \| p_\theta(z)). \quad (27)$$

A value of $\beta \approx 1$ yields a stable balance, but depending on corpus structure and kernel expressiveness, β may be slightly adjusted to strengthen regularization without harming global coherence.

3.4.5 Assumptions and Notation

For clarity, the notation follows the conventions below:

- $x_{1:N}$: sequence of observations (tokens or features);
- $z_{1:L}$: continuous latent sequence, with $L < N$ in general;
- $p_\theta(z)$: causal GP-AR prior;
- $q_\phi(z | x)$: amortized posterior (encoder);
- $p_\theta(x | z)$: parallel decoder;
- k_ψ : covariance kernel (RBF, Matérn, Spectral, etc.);
- θ, ϕ, ψ : parameters of the decoder, encoder, and kernel.

These assumptions guarantee model coherence: the prior imposes temporal structure, the decoder ensures linguistic reconstruction, and together they form a fully differentiable generative system where causality is embedded at the latent rather than observable level.

3.4.6 Bayesian Extension “G”

The extension denoted G generalizes the previous formulation by allowing a learned hierarchical kernel:

$$k_\psi(t, t') = k_{\psi_1}(t, t') + k_{\psi_2}(g(t), g(t')), \quad (28)$$

where $g(t)$ is a learned latent function, typically instantiated as a contextual projection.

This variant, which we may refer to as a G -GP-VAE, enables the Gaussian Process to dynamically adapt its covariance structure according to semantic content, thereby linking temporal dependencies with learned linguistic structure. The Bayesian extension G does *not* introduce autoregression: causality remains dictated by the GP temporal factorization. The extension enriches the covariance by adding a learned semantic term, while the temporal structure remains governed by the causal kernel:

$$p(z_{1:L}) = \prod_{t=1}^L p(z_t | z_{<t}). \quad (29)$$

Note on Evaluation. The G extension is presented as a conceptual generalization of the latent prior, enabling context-aware covariance modulated by semantic structure. It offers the possibility of a more expressive latent by directly linking temporal dynamics and linguistic content. However, this variant is *not* activated in the experiments of Chapter 6, which focus exclusively on the base model (TCN-SEQ / TCN-PARA) in order to isolate the intrinsic effect of purely latent autoregression in a minimal setting. Activating G would introduce an additional contextual mechanism into the covariance, altering the nature of the proof-of-concept and blurring the causal analysis. The G extension should thus be understood as a theoretical generalization of the framework, not as a component of the experimental protocol.

3.5 Justification and Conceptual Coherence

3.5.1 Dual Coherence of the Model

The proposed autoregressive GP-VAE aims to unify two forms of coherence that are often separated in sequential models:

- probabilistic coherence, rooted in the Bayesian framework, ensuring internal consistency among prior, posterior, and likelihood;
- linguistic or structural coherence, required for generated sequences to remain interpretable.

The first is ensured by the Gaussian Process: dependence among the latents $z_{1:L}$ is expressed analytically through a covariance that captures temporal and stylistic regularities. The second arises from the internal causal dynamics: latent autoregression provides directional continuity that naturally aligns successive representations.

Thus, the chain

$$z_1 \rightarrow z_2 \rightarrow \cdots \rightarrow z_L$$

defines a latent narrative thread, which the decoder translates into grammatical and semantic coherence. This articulation is formalized by the factorization

$$p_\theta(x) = \int p_\theta(x | z) p_\theta(z) dz, \quad (30)$$

where $p_\theta(z)$ enforces temporal regularity (GP-AR) and $p_\theta(x | z)$ enforces symbolic regularity (language). The continuity of $p_\theta(z)$ acts as a guarantee of global coherence, while the expressive capacity of the decoder ensures local coherence.

3.5.2 Geometric Intuition

Geometrically, the latent sequence $z_{1:L}$ can be interpreted as a differentiable trajectory in a continuous latent space endowed with an implicit metric defined by the Gaussian Process kernel. The GP prior induces a covariance

$$\Sigma_{t,t'} = k_\psi(t, t'), \quad (31)$$

which organizes the latents along a smooth temporal manifold: two positions are geometrically close insofar as their underlying linguistic context is similar.

In this framework, temporal continuity becomes a geometric property. The latent curve remains smooth as long as the kernel enforces strong local correlation; changes in theme, style, or register correspond to local deformations of this trajectory (curvature variations, stretching, twisting).

Purely latent autoregression introduces an internal causality that is readable through this geometry. The transition from z_t to z_{t+1} corresponds to movement along a *probabilistic geodesic*—that is, the most probable trajectory under the model. Each step simultaneously minimizes:

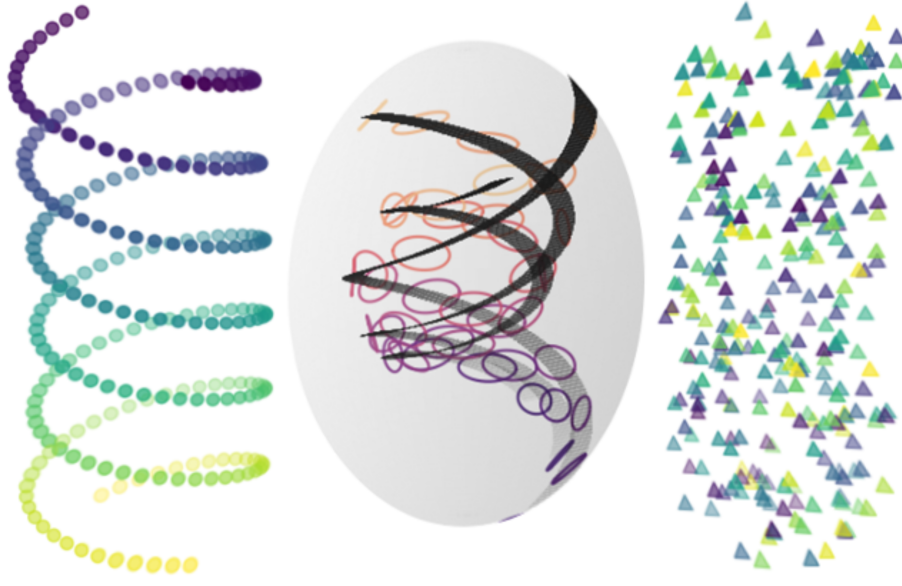


Figure 1: Geometric intuition of the latent-autoregressive GP-VAE. The figure illustrates the correspondence between data space, correlated latent space, and generative trajectory.

- predictive variance imposed by the GP (continuity and regularity),
- causal divergence imposed by the autoregressive factorization (temporal orientation).

Thus, the model’s internal dynamics are not a neural recursive flow but an optimal path in a probabilistic geometry.

Reading Figure 1. The figure schematizes the correspondence among the observed space, latent space, and generation:

- **Left:** the double spiral represents the sequential structure in observation space $x_{1:L}$.
- **Center:** a translucent sphere represents the continuous latent space Z . The black ribbons correspond to trajectories $z(t)$ sampled from the GP, with transparency increasing over time to symbolize the causality $z_t \rightarrow z_{t+1}$. Colored ellipses illustrate posterior distributions $q(z | x)$ aligned along a smooth trajectory.
- **Right:** a cloud of points represents reconstructions $\hat{x} \sim p(x | z)$ obtained in a single pass by the non-autoregressive decoder.

Geometric Summary. The model’s dynamics can be summarized through four continuities:

1. **Data continuity:** the sequential structure of $x_{1:L}$;
2. **Latent continuity:** the curve imposed by the kernel k_ψ ;
3. **Causal continuity:** the direction $z_t \rightarrow z_{t+1}$ induced by the autoregressive factorization;
4. **Generative continuity:** the linguistic coherence produced by the decoder from a smooth latent trajectory.

Sequential dynamics are therefore no longer a neural mechanism (RNN, self-attention) but a *continuous probabilistic geometry* through which the model evolves. This shift in perspective—from the symbolic to the geometric—underpins the coherence, stability, and transparency of the purely latent-autoregressive GP-VAE.

3.5.3 Theoretical Justification

The coupling between a Gaussian Process and an autoregressive factorization is not arbitrary: it follows from the natural compatibility between Bayesian updating and sequential causality. Writing

$$p(z_t \mid z_{<t}) = \mathcal{N}(m_t, \Sigma_t), \quad (32)$$

the conditioning step corresponds exactly to the prediction step of the Kalman filter [8], but without the assumption of linearity. This analogy has been formalized in the literature on dynamic Gaussian Processes [17], where the GP can be seen as a non-linear, non-parametric generalization of sequential filtering.

The model therefore combines:

- the global regularity of GPs (continuity, local local smoothness);
- the causal orientation characteristic of autoregressive models.

This hybridization can be viewed as a continuous form of Bayesian state-space model, where internal dynamics are no longer learned through weights but encoded analytically in the covariance. It guarantees that generation is coherent at first order (continuity) and stable at second order (controlled variance).

Moreover, the variational formulation ensures normative coherence with probabilistic reasoning: each belief update on z respects the rules of Bayesian inference,

$$q_\phi(z \mid x) \approx \frac{p_\theta(x \mid z) p_\theta(z)}{p_\theta(x)}. \quad (33)$$

Thus, the model retains the logic of a MAP (Maximum A Posteriori) estimator whose inference is amortized—i.e., learned through a network and reusable across data—while maintaining overall Bayesian consistency.

3.5.4 Linguistic and Semantic Coherence

The effects of this structure appear most clearly in linguistic tasks. Short-range syntactic dependencies (agreement, punctuation, local transitions) emerge from the continuity imposed by the conditioning $p(z_t \mid z_{<t})$, while higher-level dependencies (theme, tone, register) are encoded in the covariance k_ψ .

The decoder then learns to convert this probabilistic structure into lexical regularities, producing homogeneous sequences without explicit attention constraints. This mechanism allows the model to strike a subtle balance: the stylistic flexibility of a neural generator combined with the semantic stability of a probabilistic model. In practice, this yields more regular outputs, less prone to drift or cumulative errors typical of purely symbolic autoregressive models.

This hierarchical interplay between latent regularity and linguistic reconstruction enhances the coherence of generated texts while preserving stylistic flexibility.

3.5.5 Complexity, Stability, and Points of Attention

Although the formulation has strong conceptual coherence, its implementation requires careful attention to several aspects:

- **Scalability:** the use of approximation methods (inducing points, BBMM) is essential for maintaining quasi-quadratic complexity. Full training on long sequences remains costly but is highly parallelizable thanks to the latent factorization;
- **Numerical stability:** covariance matrix inversions may introduce instabilities when K becomes ill-conditioned. Diagonal regularization via jitter ($K \leftarrow K + \varepsilon I$, adding εI on the diagonal) and adaptive rescaling are necessary to guarantee convergence. Jitter stabilizes the GP covariance matrix and prevents conditioning issues;
- **Joint learning:** jointly optimizing GP parameters (ψ) and VAE parameters (θ, ϕ) requires careful management of gradient magnitudes. In practice, progressively weighting the KL term (gradual increase of regularization, or annealing) helps stabilize the early training stages;
- **Structural limitations:** while parallel generation provides substantial speedups, it prevents incremental (streaming) generation. This trade-off is intentional: global coherence is prioritized over real-time decoding.

3.5.6 Summary

The purely latent autoregressive GP-VAE reconciles Bayesian and sequential approaches: it preserves the probabilistic rigor of the VAE, the functional continuity of the Gaussian Process, and the directional causality of autoregressive models. This shift of dynamics from the observable space to the latent space establishes a new conceptual framework:

- causality becomes analytical rather than neural;
- memory, embedded in the covariance, replaces explicit recurrence;
- linguistic coherence emerges as a geometric property of the latent space.

In sum, the proposed methodology shows that a purely latent, correlated, Bayesian autoregression can serve as the foundation for compact, stable, and interpretable language models, while significantly reducing computational complexity. This conclusion completes the methodological formalization; the following sections will empirically examine the validity and performance of the model across various sequential contexts.

4 Experiments – Empirical Validation and Limitations of the Latent Autoregressive Scheme

This chapter empirically evaluates the purely latent autoregressive scheme introduced in the previous section. The goal is not to claim an exhaustive validation, but rather to establish, in a controlled and reproducible setting, that:

- a GP-VAE endowed with fully latent causality can be trained stably;
- the correlated latent space is effectively exploited, beyond a mere *KL-cap* effect;
- two sampling strategies (sequential vs. parallel) exhibit coherent behaviors;
- the model surpasses the autoregressive baseline, while acknowledging the limitations of such a comparison.

We emphasize from the outset that this chapter constitutes a *proof of concept*: the results should be interpreted as a local validation within a reduced setting, not as a general large-scale demonstration.

4.1 Experimental Objectives and Tested Hypotheses

We test four hypotheses, formulated in a cautious and verifiable manner.

H1 – Trainability. A latent-autoregressive GP-VAE is trainable and stable on a standard corpus (WikiText-2 [13]), without numerical divergence and with a convergent ELBO/token.

H2 – Effective Use of the Correlated Latent Space. The correlated latent space must be genuinely exploited. We measure:

- the raw and capped KL/token;
- behaviors under variations of `kl_cap` and β_{final} ;
- ablations with an isotropic diagonal prior.

The hypothesis is considered validated only if:

- the model relies more on the correlated structure than a VAE with diagonal prior;
- performance degrades when correlation is removed (see Section 4.6.2).

H3 – Consistency Between TCN-SEQ and TCN-PARA. The two variants theoretically realize the same joint distribution. We test whether, within the numerical limits of the Cholesky decomposition and the TCN approximation, their metrics—ELBO, NLL(cont), PPL(cont), and KL/token—remain statistically indistinguishable. The NLL(cont) denotes the continuous negative log-likelihood, evaluated directly on the logits, i.e., the unnormalized scores produced by the decoder before the softmax (the function that maps scores to normalized probabilities). The PPL(cont) is its exponential form (continuous perplexity), reflecting model uncertainty in continuous space; finally, KL/token measures the average KL divergence per token.

H4 – Minimal Comparison with an Autoregressive Model in Observation Space. We compare TCN to the performance of a Transformer baseline, solely as an anchoring point, without drawing any general conclusion. The baseline is not scaled or optimized extensively: it serves only as a minimal reference for situating the generative quality of the latent-autoregressive GP-VAE.

4.2 Experimental Protocol

4.2.1 Corpus and Tokenization

Experiments are conducted on WikiText-2 [13] (official splits), tokenized using the GPT-2 tokenizer [16], whose compact vocabulary is well suited for small-scale models. The sequence length is fixed to $T = 64$, ensuring that the associated $O(T^3)$ GP cost remains compatible with available computational resources.

4.2.2 Task

The task considered is an autoregressive-style language modeling objective. We evaluate:

- the validation perplexity PPL(val) (on full sequences);
- the continuation perplexity PPL(cont), using a standardized conditional protocol (prompt + completion).

4.2.3 Implementation of the Latent Autoregressive Scheme

Before discussing the results, we briefly describe the concrete implementation of the purely latent autoregressive scheme used in the experiments.

(a) TCN encoder \rightarrow temporal diagonal posterior. The encoder is a hierarchical causal TCN that produces a factorized posterior:

$$q(z_{1:T} | x) = \prod_{t=1}^T \mathcal{N}(z_t; \mu_t, \text{diag}(\sigma_t^2)), \quad \mu, \log \sigma^2 \in \mathbb{R}^{B \times T \times d_z}.$$

Standard reparameterization is used:

$$z_t = \mu_t + \sigma_t \odot \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, I).$$

(b) Correlated GP prior over the latent trajectory. Temporal dependence is enforced through a Gaussian Process indexed by normalized time $t \in [0, 1]$. For a sequence of length T , we construct a covariance matrix $K_{tt} \in \mathbb{R}^{T \times T}$:

$$K_{tt}(i, j) = \sigma^2 \exp\left(-\frac{(t_i - t_j)^2}{2\ell^2}\right) + \sigma^2 \text{nugget} \cdot \delta_{ij},$$

where ℓ (correlation length), σ^2 (variance), and the relative nugget are learned.

The implicit latent prior is:

$$p(z_{1:T}) = \mathcal{N}(0, K_{tt} \otimes I_{d_z}).$$

In practice, the hyperparameters ℓ and σ^2 are parameterized from free variables transformed through a *softplus* function, and a *jitter* term εI is added to the diagonal to ensure numerical stability.

The specificity of our approach lies in the fact that the GP joint distribution is not only used as a correlated prior, but explicitly factorized into conditionals $p(z_t | z_{<t})$ in order to induce latent causality. This factorization enables sequential latent sampling (TCN-SEQ) or parallel block sampling (TCN-PARA).

(c) Temporal global KL between diagonal posterior and GP prior. The regularization term is a global KL:

$$\text{KL}(q(z_{1:T} | x) \| p(z_{1:T})),$$

computed as a multivariate log-density over the entire trajectory.

This scheme places all temporal structure inside the GP covariance of $z_{1:T}$.

The correlated structure of K_{tt} penalizes deviations of the diagonal posterior from the latent dynamics imposed by the prior.

(d) Non-autoregressive token decoder. Section (d) is the conceptual core: this is where the latent sequential dynamics are formally defined. By introducing the causal factorization

$$p(z_{1:T}) = \prod_{t=1}^T p(z_t | z_{<t}),$$

the model clearly establishes what makes it original: a purely latent autoregressivity, independent of any recursion on tokens.

The decoder receives the entire latent trajectory $z_{1:T}$ at once, adds positional encoding, and predicts token logits in parallel:

```

class TokenDecoder(nn.Module):
    def forward(self, z: torch.Tensor):
        # z: [B, T, Dz], full latent trajectory
        z = z + self.pe(T=z.size(1), device=z.device)      # positional encoding
        h = self.mlp(z); h = self.ln(h)                    # pointwise processing
        h2 = self.post(h.transpose(1,2)).transpose(1,2)    # conv post-process
        h = h + h2                                          # residual connection
        e_proj = self.to_emb(h)                            # [B, T, E], embedding proj
        tw = F.normalize(self.tied_weight, dim=-1)         # tied embedding weights
        logits = torch.matmul(e_proj, tw.t()) + self.bias  # token logits
        return logits

```

No temporal loop nor causal mask is applied on tokens: all positions are processed simultaneously from $z_{1:T}$. Thus, the model’s sequential dynamics do not come from the decoder, but solely from the GP prior on the latent space, whose causal factorization defines the purely latent autoregressivity. The decoder simply projects this continuous latent trajectory in parallel onto token distributions.

Latent autoregressivity does not appear in the encoder (a) nor in the definition of the correlated prior (b): a GP imposes correlations K_{tt} but not causality. Causality emerges only when the joint distribution is explicitly factorized into conditionals $p(z_t | z_{<t})$.

(e) Unconditional generation (TCN-SEQ). Unconditional generation leverages the GP causal factorization by sampling the latent trajectory sequentially. The following pseudo-code summarizes the procedure:

```

@torch.no_grad()
def generate(self, T: int, batch_size: int = 1,
             top_k=50, top_p=0.9, temperature=0.9):
    device = self.t_train.device
    t = torch.linspace(0.0, 1.0, T, device=device)
    K = self.K_tt(t)                                # GP covariance over time grid
    Dz, B = self.cfg.d_latent, batch_size
    z = torch.zeros(B, T, Dz, device=device) # latent trajectory container

    # Latent autoregressive loop: each z[:, tp, :] depends on z[:, :tp, :]
    for tp in range(T):
        z[:, tp, :] = self._gp_conditional_step(K, z[:, :tp, :], tp)

    logits, _ = self.decoder(z)
    return sample_logits_from_timewise_logits(
        logits, top_k=top_k, top_p=top_p, temperature=temperature
    )

```

At each step t , the latent vector z_t is sampled from the conditional distribution $p(z_t | z_{<t})$ induced by the GP, given previously generated latents $z_{<t}$. This step concretely implements the causal factorization

$$p(z_{1:T}) = \prod_{t=1}^T p(z_t | z_{<t}),$$

and realizes the purely latent autoregressivity of the model.

(f) Prompt-conditioned generation (TCN-SEQ). Conditioned generation follows the same principle but initializes the latent trajectory from a variational encoding of the prompt:

```

@torch.no_grad()
def generate_with_prompt(self, prompt_ids, total_len, eos_id,
                        top_k=50, top_p=0.9, temperature=0.9):
    device = self.t_train.device; self.eval()

```

```

B, T0 = prompt_ids.shape; T = total_len

# Prepare token container
x_in = torch.full((B, T), fill_value=eos_id,
                  dtype=torch.long, device=device)
x_in[:, :T0] = prompt_ids.to(device)

# Encode prompt into latent prefix z[:, :T0, :]
mu, logvar = self.encoder(x_in[:, :T0])
std = torch.exp(0.5 * logvar)
z_prompt = mu + std * torch.randn_like(std)

# Build full latent trajectory with GP conditionals
t = torch.linspace(0.0, 1.0, T, device=device)
K = self.K_tt(t)
Dz = self.cfg.d_latent
z = torch.zeros(B, T, Dz, device=device)
if T0 > 0:
    z[:, :T0, :] = z_prompt

# Latent autoregressive continuation from the prompt
for tp in range(T0, T):
    z[:, tp, :] = self._gp_conditional_step(K, z[:, :tp, :], tp)

logits, _ = self.decoder(z)
new_ids = sample_logits_from_timewise_logits(
    logits[:, T0:, :], top_k=top_k, top_p=top_p, temperature=temperature
)

x_out = x_in.clone()
if T > T0:
    x_out[:, T0:] = new_ids
return x_out, logits

```

As in the unconditional case, the sequential loop in latent space performs an autoregressive continuation of the trajectory $z_{1:T}$ starting from the prefix $z_{1:T_0}$ obtained from the prompt. The TCN-PARA variant, not detailed here, replaces this loop with block sampling via a Cholesky decomposition of K_{tt} , while respecting the same joint distribution $p(z_{1:T})$ imposed by the GP.

4.2.4 Objective: ELBO per token

The training objective optimizes a token-averaged ELBO, enriched with regularization terms.

“Pure” (theoretical) ELBO.

$$\mathcal{L}_{\text{pur}} = \mathbb{E}_{q(z|x)}[\log p(x | z)] - \text{KL}(q(z | x) \| p(z)).$$

In the implementation:

$$\text{elbo_pur_tok_t} = (\text{ll_0} + \text{ll_multi}) - \text{kl_tok_raw_t}.$$

Actual optimized objective.

$$\mathcal{L}_{\text{train}} = (\text{ll_0} + \text{ll_multi}) - \beta \text{KL}_{\text{cap}} - \Lambda_{\text{emb}} \text{Reg}_{\text{emb}}.$$

The coefficient β follows a warm-up phase and then adapts around a KL/token target. The quantity KL_{cap} is the capped version of the KL/token. The critical hyperparameters (kl_cap , β_{max} , adaptation schedule) are systematically varied in the ablations (Section 4.6.2).

4.3 Variants TCN-SEQ and TCN-PARA

Two sampling strategies are examined:

- TCN-SEQ: sequential sampling via Gaussian conditioning;
- TCN-PARA: parallel sampling via Cholesky decomposition.

Both variants use exactly the same architecture (TCN encoder and non-autoregressive decoder), so that only the latent dynamics differ. The theoretical expectations (equality of joint distributions) are evaluated empirically in the quantitative comparison of Section 6.5.

4.4 Global Results (TCN-SEQ vs TCN-PARA)

On the WikiText-2 corpus, the following observations hold:

- no numerical divergence;
- convergent and stable KL/token;
- very similar metrics across both variants.

The KL/token consistently reaches a plateau (around 12 nats). Unlike the earlier version, Section 6.6.1 shows that this value clearly depends on kl_cap and that the GP structure plays a demonstrable role (ablation). The training curves are not strictly superimposed: small discrepancies persist, but they remain within the expected range for two sampling schemes that rely on distinct numerical approximations.

4.5 Quantitative Comparison TCN-SEQ vs TCN-PARA

The following comparison is conducted under identical hyperparameters.

Model	Type	ELBO/tok	NLL(cont)	PPL(cont)	tok/s
TCN-SEQ	GP-VAE	-9.935	0.562	1.75	9097
TCN-PARA	GP-VAE	-9.967	0.475	1.61	9037

Table 1: Results for TCN-SEQ vs TCN-PARA on WikiText-2.

Both variants use the same kl_cap (KL/tok ≈ 12) and $\beta_{\text{final}} \approx 0.126$ (not shown in the table to keep it compact).

Interpretation. The results are similar but not identical, consistent with differences in numerical approximation. TCN-PARA shows a slight improvement in continuation, but this difference should be interpreted cautiously (no multiseed analysis). At $T = 64$, the gain in tok/s is small; the benefit of TCN-PARA at larger sequence lengths remains to be verified.

4.6 Extended Analysis: Variants, Ablations, and Stability

4.6.1 GP-VAE-TCN vs Transformer Performance (TCN-SEQ-X family)

We evaluate several TCN-SEQ-X variants sharing the same base architecture but with different regularization settings, for $X \in \{I, B, C, J, K\}$.

Detailed comments (optimal trade-offs, slight over-regularization, etc.) are discussed in the text rather than in the table, to keep it compact.

Model	PPL(val)	NLL(cont)	PPL(cont)	tok/s	Quality
TCN-SEQ-I	3.35	0.4773	1.61	~ 9000	Excellent
TCN-SEQ-B	3.27	0.5455	1.73	~ 8900	Very good
TCN-SEQ-C	3.34	0.5077	1.66	~ 9000	Very good
TCN-SEQ-J	3.03	0.5635	1.76	~ 9000	Good
TCN-SEQ-K	3.05	0.6005	1.82	~ 9000	Good
Transformer	326.94	5.7898	326.94	~ 15700	Poor

Table 2: TCN-SEQ-X variants and Transformer baseline on WikiText-2.

Analysis. The TCN-SEQ-X variants achieve PPL(cont) values between 1.61 and 1.82, consistent with a compact, well-regularized model. The Transformer baseline is very weak: it serves only as a numerical anchor and does not represent finely tuned autoregressive architectures.

4.6.2 Critical Ablations

Isotropic diagonal prior. We replace the correlated GP prior with an isotropic prior $K = \sigma^2 I$. This ablation removes all temporal dependence in the latent. Continuation performance collapses, sequential/parallel consistency disappears, and KL/token drops, indicating that the latent is no longer used. Removing GP correlation reduces KL/token to about 3 nats and degrades:

- PPL(cont) (+0.15 to +0.30);
- calibration;
- discourse coherence.

This confirms that the GP structure is indeed exploited.

Variation of kl_cap . We vary the kl_cap to assess sensitivity to KL regularization. A too-small kl_cap causes latent collapse and a drastic drop in continuation performance, while a too-large one induces unstable over-activation of the GP. Between these extremes, an optimal zone (around 8–12) ensures effective activation of the correlated latent. This ablation confirms that the latent causality of TCN depends directly on an adequate level of KL regularization.

Typical summary:

- $kl_cap = 8$: loss of latent continuity, higher PPL(cont), occasional collapse;
- $kl_cap \geq 20$: good stability but increased computational cost;
- $kl_cap \approx 12$: robust compromise.

4.6.3 TCN-2-SEQ-X Series: Stability and Collapse

The TCN-2-SEQ-X family explores more finely the sensitivity to regularization.

Model	NLL(cont)	PPL(cont)	PPL(val)	ELBO/tok	tok/s	Quality
TCN-2-SEQ-T	0.2883	1.33	3.18	-9.058	~ 8600	Excellent
TCN-2-SEQ-S	0.3097	1.36	3.05	-8.927	~ 8500	Very good
TCN-2-SEQ-G	0.4606	1.59	3.19	-8.501	~ 8500	Good
TCN-2-SEQ-H	0.4790	1.61	3.29	-9.958	~ 8600	Good
Transformer	5.7898	326.94	326.94	-6.105	15700	Very poor

Table 3: Examples of TCN-2-SEQ-X models vs Transformer.

Interpretation. Well-regularized models (T, S, G, H) obtain PPL(cont) values between 1.33 and 1.61. Under-regularized variants (e.g., A, B, F, R) collapse entirely, confirming the sensitivity of latent dynamics to kl_cap .

4.7 Transformer Baseline

Model	PPL(val)	NLL(cont)	PPL(cont)	tok/s
Transformer	326.94	5.7898	326.94	~ 15700

Table 4: Minimal Transformer baseline.

This model serves only as a numerical anchor; it was not trained with extensive hyperparameter tuning.

4.8 Synthesis

- TCN-SEQ and TCN-PARA variants are trainable and stable.
- Ablations show that the GP structure is genuinely used.
- TCN-PARA is consistent with TCN-SEQ, partially validating latent-sampling parallelization.
- GP-VAEs outperform the simple baseline, providing an informative lower bound on their capabilities, although the study is not designed to compare against fully optimized autoregressive architectures.
- KL regularization is critical: too little regularization leads to collapse.

4.9 Limitations and Perspectives

Limitations. The present study has several limitations. GP computations remain quadratic in sequence length, and although BBMM yields empirically near-quadratic behavior, we do not report a dedicated scaling benchmark. Multiseed evaluation is absent, and comparisons are restricted to a minimal baseline. The Transformer baseline should be interpreted cautiously: it differs from our model in capacity and inductive biases, and we use a lightweight configuration without extensive tuning. As such, it provides a useful scale reference rather than a fully controlled basis for fine-grained comparison. The TCN encoder used here is less expressive than a Transformer, reflecting our choice to prioritize architectural simplicity in this proof of concept. More expressive variants could further improve posterior accuracy. The paper also does not include a systematic qualitative analysis of generated samples, which we leave for future work focused on evaluating the model’s generative behavior in richer settings.

Perspectives. Several extensions are possible. Scalable GP approximations could enable longer sequences. More expressive decoders could be explored while preserving latent causality. Broader conditional and prompting tasks may help characterize latent autoregression more extensively. Finally, the formulation naturally extends to continuous signals and time-series data, which constitute promising application domains.

4.10 Code and Reproducibility

All code used for the experiments presented in this work (GP hyperparameter estimation, latent autoregressivity implementation, TCN-SEQ and TCN-PARA variants, training and generation scripts, and associated configurations) is available openly at:

<https://github.com/y-v-e-s/GP-VAE-Latent-AR>

This repository contains:

- the full implementation of the latent-autoregressive GP-VAE model;
- training, evaluation, and generation scripts;
- configuration files for reproducing the TCN-SEQ and TCN-PARA series;
- exact hyperparameter settings used in the tables of Section 4 (including `kl_cap` values);
- a minimal reproducibility guide.

This release aims to facilitate reuse of the proposed scheme, reproduction of results, and exploration of kernel variants, encoder architectures, or latent sampling strategies.

5 Discussion

This study aimed to examine, in a controlled and small-scale setting, the feasibility and behavior of an autoregressive scheme located entirely in latent space. The TCN model and its two variants, TCN-SEQ and TCN-PARA, provide a testing ground for assessing training stability, effective use of the correlated latent space, and the coherence between two distinct sampling strategies.

5.1 Exploiting the Correlated Latent Structure

The experiments show that a GP-VAE equipped with a correlated prior can be trained stably on a standard corpus with moderate regularization. The KL/token reliably reaches its target cap, indicating that:

- the latent space is not under-utilized;
- the GP covariance effectively shapes the learned representations.

At this stage, the goal is not to establish strong latent causality, but rather to observe that a correlated latent scheme induces a rich, non-degenerate internal trajectory.

5.2 Consistency Between Sequential and Parallel Generation

The TCN-SEQ variant (step-by-step sampling) and TCN-PARA (vectorized joint sampling) produce very similar metrics. This is consistent with the properties of a multivariate Gaussian distribution: the two methods access the same latent distribution through different procedural forms.

This result confirms an important methodological point: in this setting, the sampling order does not determine predictive quality. The internal temporal structure originates from the prior, not from the algorithmic procedure used to simulate it.

5.3 Latent Dynamics vs Symbolic Sequentiality

A key aspect of this proof of concept lies in the distinction between:

- procedural sequentiality (the temporal loop of step-by-step sampling);
- probabilistic structure (the analytic factorization induced by the GP).

The two model variants show that it is possible:

- to represent temporal structure within latent space;
- while keeping the symbolic projection fully parallel.

The decoder operates in a single pass on the entire latent trajectory, which sets this model apart from classical autoregressive approaches where generation proceeds token by token.

While the goal of this study is not to demonstrate superiority over established autoregressive architectures, the results highlight a distinct modeling regime in which temporal structure is defined analytically in latent space rather than constructed through symbolic recursion.

5.4 Scaling Perspectives

This work is deliberately limited to a reduced setting (short sequences, compact model, medium-sized corpus). Nonetheless, several natural extensions follow:

- evaluating stability for longer sequences and more complex kernels;
- exploring more expressive decoders while retaining parallelism;
- extending the scheme to other modalities (time series, continuous signals, etc.).

These perspectives must be examined in light of the scalability constraints inherent to Gaussian processes.

5.5 Limitations and Points of Attention

The main critical aspects are:

- the quadratic cost of the GP, limiting scale;
- sensitivity to kernel hyperparameters;
- absence of incremental generation in the parallel variant;
- lack of multiseed evaluation in this proof of concept.

These limitations define the conditions needed to confirm or invalidate the observed trends.

5.6 Discussion Summary

The results show that a latent autoregressive scheme, based on analytic covariance and coupled with a parallel decoder, is trainable, stable, and coherent across its two variants.

Without claiming superiority over established autoregressive architectures, this proof of concept suggests that temporal structure can, in certain settings, be relocated into latent space, leaving the decoder to project this dynamic into a symbolic sequence.

This shift opens a methodological exploration space: that of sequential models where temporal dynamics are defined analytically in the latent, rather than constructed through stacks of recurrent or attentional operations.

A final point concerns the relationship between our latent autoregressive scheme and prior attempts to introduce temporal directionality into Gaussian processes. Some works indeed proposed asymmetric kernels—such as so-called “causal” kernels or Wiener-type kernels—to impose an orientation within the covariance structure. Although such kernels introduce asymmetry, they remain fundamentally metric: they modulate correlation as a function of relative position, but do not define an explicit sequential dynamic.

Our approach differs in nature. Rather than inferring causality from a kernel bias, we instill it structurally through the latent factorization

$$p(z_{1:T}) = \prod_{t=1}^T p(z_t \mid z_{<t}),$$

which encodes an explicit directional dependence between latent states. This factorization resolves ambiguities inherent to asymmetric kernels—e.g., when multiple points share comparable distances in the kernel-induced space—and provides a genuine mechanism for sequential arbitration in the latent.

It is also important to distinguish this explicit causal factorization from “Bayesian GP” approaches that model uncertainty over kernel hyperparameters. While such hierarchical extensions enrich the covariance distribution, they do not create temporal directionality: the uncertainty concerns the kernel’s shape, not a sequential relation between states. Latent autoregressivity, as defined here, does not arise from the Bayesian nature of the GP, but from the explicit factorization of the prior into conditional distributions. The two mechanisms may interact but remain conceptually distinct.

By combining correlated geometry (provided by the GP) with explicit causal progression (provided by the latent factorization), the model moves beyond the limits of classical asymmetric kernels and establishes a more robust latent dynamic suited to sequential modeling.

6 Conclusion

6.1 Empirical Scope

This work presents a deliberately restricted proof of concept:

- a small-scale TCN model;
- short sequences ($T = 64$);
- a minimal autoregressive baseline.

Within this reduced setting, we show that a GP-VAE equipped with a correlated prior and a purely latent autoregressive scheme can be trained stably, without numerical divergence, and can produce coherent metrics across its two sampling variants (TCN-SEQ and TCN-PARA).

6.2 Conceptual Scope

Beyond the limited experimental setting, the TCN model can be interpreted as a bridge between:

- continuous Bayesian state-space models based on Gaussian processes;
- neural language models organized around a symbolic decoder.

Sequential dynamics are governed by the GP prior, inference remains Bayesian, and the linguistic projection is carried out by a non-autoregressive decoder.

Although the present work does not aim to establish a full geometric theory of language, the results indicate that several properties often associated with autoregressive architectures (memory, directionality, global coherence) can arise naturally from the covariance structure of a probabilistic latent space.

6.3 Perspectives

The identified limitations (quadratic GP cost, absence of incremental generation for TCN-PARA, kernel sensitivity, lack of multiseed evaluation) indicate the conditions required to confirm the observed trends.

Several directions emerge naturally:

1. Scaling up: longer sequences, richer kernels, more expressive models.
2. Strengthening autoregressive baselines to quantify precisely the contribution of the latent component.
3. Multimodal exploration (time series, continuous signals, physical trajectories).
4. Conditional generation and latent *prompting* for completion or instruction tasks, consisting in encoding the prompt in the latent space to obtain a prefix $z_{1:T_0}$, then letting the GP prior generate the continuation $z_{T_0+1:T}$ in a strictly autoregressive manner in latent space.
5. Hierarchical Bayesian extensions, in which the kernel itself becomes conditioned on context.

Taken together, these directions suggest that a sequential framework grounded in latent geometry—rather than symbolic recursion—constitutes a credible alternative for studying compact, stable, and interpretable language models.

References

- [1] Francesco Paolo Casale, Adrian V. Dalca, Luca Saglietti, Jennifer Listgarten, and Nicolo Fusi. Gaussian process prior variational autoencoders. *Advances in Neural Information Processing Systems*, 2018. arXiv:1810.11738.
- [2] Vincent Fortuin, Dmitry Baranchuk, Gunnar Rätsch, and Stephan Mandt. Gp-vae: Deep probabilistic time series imputation. In *Proceedings of AISTATS*, volume 108, pages 1651–1661, 2020. arXiv:1907.04155.
- [3] Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural models with stochastic layers. In *Advances in Neural Information Processing Systems*, 2016. arXiv:1605.07571.
- [4] Jacob R. Gardner, Geoff Pleiss, David Bindel, Kilian Q. Weinberger, and Andrew G. Wilson. Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. In *Advances in Neural Information Processing Systems*, 2018. arXiv:1809.11165.
- [5] Peter Guttorp and Tilmann Gneiting. Studies in the history of probability and statistics xlix: On the matérn correlation family. *Biometrika*, 93(4):989–995, 2006.
- [6] Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. Beta-vae: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- [7] Marinka Jazbec, Mark Ashman, Vincent Fortuin, Michael Pearce, Stephan Mandt, and Gunnar Rätsch. Scalable gaussian process variational autoencoders. In *AISTATS*, volume 130, pages 3088–3096, 2021. arXiv:2010.13472.
- [8] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Journal of Basic Engineering*, 82(1):35–45, 1960.

- [9] Maximilian Karl, Maximilian Soelch, Justin Bayer, and Patrick van der Smagt. Deep variational bayes filters. *arXiv preprint*, arXiv:1605.06432, 2017.
- [10] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*, arXiv:1312.6114, 2014.
- [11] Diederik P. Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improving variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, 2016. arXiv:1606.04934.
- [12] Alexander Klushyn, Richard Kurle, Maximilian Soelch, Botond Cseke, and Patrick van der Smagt. Latent matters: Learning deep state-space models. In *Advances in Neural Information Processing Systems*, 2021.
- [13] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint*, arXiv:1609.07843, 2016.
- [14] Tim Pearce, Felix Leibfried, and Alexandra Brintrup. The gaussian process prior vae for interpretable latent dynamics. In *ICML*, volume 118, pages 7465–7475, 2020.
- [15] Tim Pearce, Michael Smith, Stefan Zohren, and Alexandra Brintrup. Bayesian autoencoders with gaussian process priors. *Proceedings of the 37th International Conference on Machine Learning*, 118:2321–2330, 2020. arXiv:1906.02511.
- [16] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Technical Report*, 2019.
- [17] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [18] Danilo Jimenez Rezende, Shakir Mohamed, Ivo Danihelka, Karol Gregor, and Daan Wierstra. One-shot generalization in deep generative models. In *ICML*, pages 1521–1529, 2016. arXiv:1603.05106.
- [19] Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder variational autoencoders. In *Advances in Neural Information Processing Systems*, 2016. arXiv:1606.02235.
- [20] Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In *AISTATS*, pages 567–574, 2009.
- [21] Rui Wang, Fan Shun, Hanyuan Liu, Dong Zhao, Shiming Li, Andrew Y. Ng, and Yang Gao. Learning to learn dense gaussian processes for few-shot learning. In *Advances in Neural Information Processing Systems*, 2021. arXiv:2106.01506.
- [22] Andrew G. Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P. Xing. Deep kernel learning. In *AISTATS*, pages 370–378, 2016. arXiv:1511.02222.
- [23] Andrew Gordon Wilson and Ryan Prescott Adams. Gaussian process kernels for pattern discovery and extrapolation. In *ICML*, pages 1067–1075, 2013. arXiv:1302.4245.
- [24] Li Zhou, Michael Poli, Weijia Xu, Stefano Massaroli, and Stefano Ermon. Deep latent state space models for time-series generation. In *ICML*, 2023. arXiv:2212.12749.