

Latent Optimal Paths by Gumbel Propagation for Variational Bayesian Dynamic Programming

Xinlei Niu¹, Christian Walder², Jing Zhang¹, Charles Patrick Martin¹

¹Australian National University

²Google DeepMind

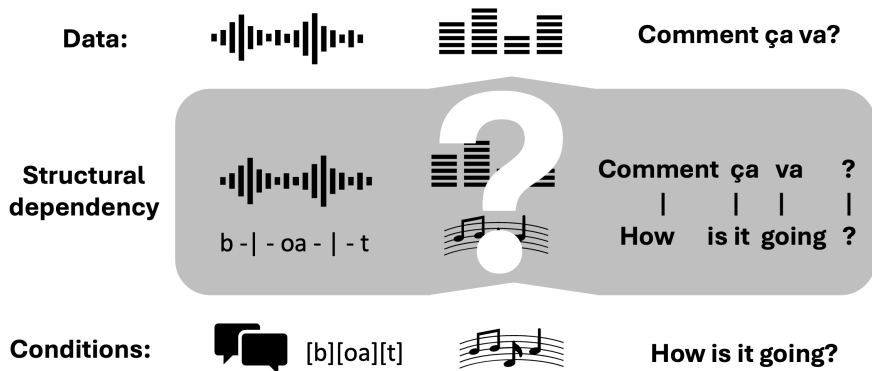


**Australian
National
University**



Google DeepMind

Optimal Paths in Generative Tasks



Related Work

- Traditional optimal path problem
 - ① Loose the sparsity of optimal paths with differentiable dynamic programming (DP) (Verdu & Poor, 1987; Kim et al., 2017; Cuturi & Blondel, 2017).
 - ② Involve a convex optimization problem (Amos & Kolter, 2017; Djolonga & Krause, 2017).
 - ③ Maintain the sparsity but smoothed with a convex regularizer (Mensch & Blondel, 2018)
- Downstream generative tasks with unobserved structural dependency
 - ① Multiple training strategies, strongly relies on external components (Ren et al., 2020; Liu et al., 2022).
 - ② Have a strong assumption about the model structure (Rabiner, 1989; Petrov & Klein, 2007).
 - ③ Cannot ensure train and test consistency (Kim et al., 2020) or soft approximation (Shen et al., 2018; Lu et al., 2021).

Our Motivation: Stochastic Optimal Paths

Definition

For every possible path $\mathbf{y} \in \mathcal{Y}$, denote by

$$\mathcal{D}(\mathcal{R}, \mathbf{W}, \alpha) \quad (1)$$

Gibbs distribution over a path $\mathbf{y} \in \mathcal{Y}(1, N)$ with probability mass function

$$\mathcal{D}(\mathbf{y}|\mathcal{R}, \mathbf{W}, \alpha) = \frac{\exp(\alpha \|\mathbf{y}\|_{\mathbf{W}})}{\sum_{\hat{\mathbf{y}} \in \mathcal{Y}(1, N)} \exp(\alpha \|\hat{\mathbf{y}}\|_{\mathbf{W}})}. \quad (2)$$

- 1 A **unified** method to obtain structured **sparse** optimal paths.
- 2 Provide **tractable** closing forms for all ingredients of variational Bayesian inference to capture stochastic optimal paths.
- 3 Allow **gradient optimization** for probabilistic generative models.

Definition of a Directed Acyclic Graph

- Denote $\mathcal{R} = (\mathcal{V}, \mathcal{E})$ be a directed acyclic graph (DAG) with nodes \mathcal{V} and edges \mathcal{E} .
- The nodes are numbered in topological order, s.t., $\mathcal{V} = (1, 2, \dots, N)$ and $u < v$ for all $(u, v) \in \mathcal{E}$.
- Assume 1 is the only node without parents and N is the only node without children.
- Denote the edge weights matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$ with $w_{i,j} = -\infty$ for all $(u, v) \notin \mathcal{E}$.
- Let $\mathcal{Y}(1, v)$ be the set of all paths from nodes $(1, \dots, v)$.
- Each path $\mathbf{y} = (y_1, \dots, y_{|\mathbf{y}|})$ has a score by summing edge weights along the path, defined as $\|\mathbf{y}\|_{\mathbf{W}} = \sum_{(u,v) \in \mathbf{y}} w_{u,v}$.

Gumbel Random Variable

Let $\mathcal{G}(\mu)$ denote the unit scale Gumbel random variable with location parameter μ and probability density function.

$$\mathcal{G}(x|\mu) = \exp(-(x - \mu) - \exp(x - \mu)). \quad (3)$$

- Shifting: Let $X \sim \mathcal{G}(\mu)$.

$$X + \text{const.} \sim \mathcal{G}(\mu + \text{const.}). \quad (4)$$

- Max: Let $X_i \sim \mathcal{G}(\mu_i)$ for all $i \in \{1, 2, \dots, m\}$.

$$\max(\{X_1, X_2, \dots, X_m\}) \sim \mathcal{G}(\log \sum_{i=1}^m \exp(\mu_i)). \quad (5)$$

- Argmax:

$$p(k = \operatorname{argmax}_{i \in \{1, 2, \dots, m\}} X_i) = \frac{\exp(\mu_k)}{\sum_{i=1}^m \exp(\mu_i)}. \quad (6)$$

Gumbel Propagation

Gumbel propagation offers an equivalent formulation of the definition that lends itself to dynamic programming by Gumbel max and shift properties.

Lemma

Let

$$Y = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(1, N)} \{\Omega_{\mathbf{y}}\}, \quad (7)$$

where for all $\mathbf{y} \in \mathcal{Y}(1, N)$,

$$\Omega_{\mathbf{y}} = \alpha \|\mathbf{y}\|_{\mathbf{W}} + G_{\mathbf{y}} \quad (8)$$

$$G_{\mathbf{y}} \sim \mathcal{G}(0). \quad (9)$$

Then the probability of $Y = \mathbf{y}$ is given by (2).

Gumbel Propagation (Cont.)

Let the definitions of $\Omega_{\mathbf{y}}$ and $G_{\mathbf{y}}$ extend to all $\mathbf{y} \in \bigcup_{u=1}^N \mathcal{Y}(1, u)$, which is the set of all partial paths. We define for each node $v \in \mathcal{V}$ the real-valued random variable

$$Q_v = \max_{\mathbf{y} \in \mathcal{Y}(1, v)} \{\Omega_{\mathbf{y}}\}. \quad (10)$$

Lemma

The Q_v are Gumbel distributed with

$$Q_v \sim \mathcal{G}(\mu_v), \quad (11)$$

where

$$\mu_v = \log \sum_{\mathbf{y} \in \mathcal{Y}(1, v)} \exp(\alpha \|\mathbf{y}\|_{\mathbf{W}}). \quad (12)$$

Gumbel Propagation (Cont.)

We now state our main result:

Lemma

The location parameters μ_v satisfy the recursion

$$\mu_1 = 0 \tag{13}$$

$$\mu_v = \log \sum_{u \in \mathcal{P}(v)} \exp(\mu_u + \alpha w_{u,v}). \tag{14}$$

for all $v \in \{2, 3, \dots, N\}$.

Gumbel Propagation (Cont.)

Lemma

Let paths $\mathbf{y} = (y_1, y_2, \dots, y_{|\mathbf{y}|})$ denote the component of the random variable Y defined in (7), given that $Y = (y_1, y_2, \dots, y_{|\mathbf{y}|})$. The probability of the transition $v \rightarrow u$ is

$$\pi_{u,v} \equiv p(y_{i-1} = u | y_i = v, u \in \mathcal{P}(v)) \quad (15)$$

$$= \frac{\exp(\mu_u + \alpha w_{u,v})}{\exp(\mu_v)}, \quad (16)$$

for all $i \in \{2, 3, \dots, N\}$.

Bayesian Dynamic Programming

- Sampling

Corollary

Paths $\mathbf{y} \sim \mathcal{D}(\mathcal{R}, \mathbf{W}, \alpha)$ may be sampled (in reverse) by

- 1 *Initializing $v = N$,*
- 2 *sampling $u \in \mathcal{P}(v)$ with probability $\pi_{u,v}$,*
- 3 *setting $v \leftarrow u$,*
- 4 *if $v = 1$ then stop, otherwise return to step 2.*

- Likelihood
- KL Divergence within the distribution family of $\mathcal{D}(\mathcal{R}, \cdot, \alpha)$.

Bayesian Dynamic Programming (Cont.)

- Sampling
- Likelihood

Corollary

The path probability may be written

$$\mathcal{D}(\mathbf{y}|\mathcal{R}, \mathbf{W}, \alpha) = \prod_{(u,v) \in \mathbf{y}} \pi_{u,v}. \quad (17)$$

- KL Divergence within the distribution family of $\mathcal{D}(\mathcal{R}, \cdot, \alpha)$.

Bayesian Dynamic Programming (Cont.)

- Sampling
- Likelihood
- KL Divergence within the distribution family of $\mathcal{D}(\mathcal{R}, \cdot, \alpha)$.

Lemma

The KL divergence within the family $\mathcal{D}(\mathcal{R}, \cdot, \alpha)$ is

$$\begin{aligned} \mathcal{D}_{\text{KL}} \left[\mathcal{D}(\mathcal{R}, \mathbf{W}, \alpha) \parallel \mathcal{D}(\mathcal{R}, \mathbf{W}^{(r)}, \alpha) \right] \\ = \mu_N^{(r)} - \mu_N + \alpha \sum_{(u,v) \in \mathcal{E}} \omega_{u,v} (w_{u,v} - w_{u,v}^{(r)}), \end{aligned} \quad (18)$$

where $\omega_{u,v}$ is the marginal probability of edge (u, v) on $\mathcal{D}(\mathcal{R}, \mathbf{W}, \alpha)$, $\mu_N^{(r)}$ is similar to μ_N but defined in terms of $\mathbf{W}^{(r)}$ rather than \mathbf{W} .

BDP-VAE

Find unobserved dependencies between condition and data in a latent space, thereby, helps a better reconstruction ¹.

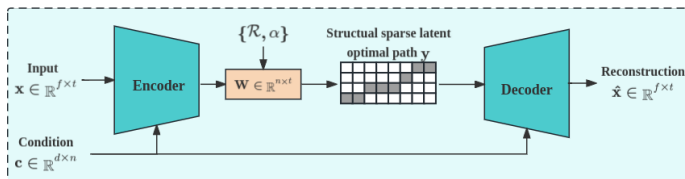


Figure: A pipeline of BDP-VAE, where we wish to find an unobserved hard structural relationship between \mathbf{x} and \mathbf{c} in the latent space of VAEs denoted as \mathbf{y} .

¹For soft optimal paths, please refer to our discussion in [Appendix.D](#)

Application on Downstream Tasks

① End-to-end Text-to-Speech (TTS)

- ▶ We adapt a non end-to-end TTS model (FastSpeech2) to BDP-VAE framework, thereby performs an end-to-end TTS model (i.e., BDPVAE-TTS).
- ▶ BDPVAE-TTS outperforms than the baseline method with an end-to-end pipeline.
- ▶ BDPVAE-TTS ensures train and test consistency on finding the monotonic alignment path.

Table: Mel Cepstral Distortion (MCD) and Real-Time Factor (RTF) compared with other TTS models.

Model	Training	Align.(Train)	Align. (Infer)	MCD	RTF
FastSpeech2 (Baseline)	Non end-to-end	Discrete	Continuous	9.96 ± 1.01	3.87×10^{-4}
Tacotron2	End-to-end	Continuous	Continuous	11.39 ± 1.95	6.07×10^{-4}
VAENAR-TTS	End-to-end	Continuous	Continuous	8.18 ± 0.87	1.10×10^{-4}
Glow-TTS	End-to-end	Discrete	Continuous	8.58 ± 0.89	2.87×10^{-4}
BDPVAE-TTS (ours)	End-to-end	Discrete	Discrete	8.49 ± 0.96	3.00×10^{-4}

Application on Downstream Tasks

- ① End-to-end Text-to-Speech (TTS)
- ② **End-to-end Singing Voice Synthesis**

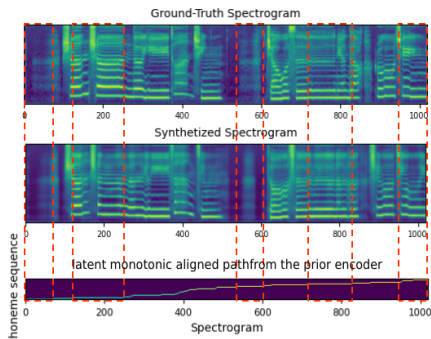
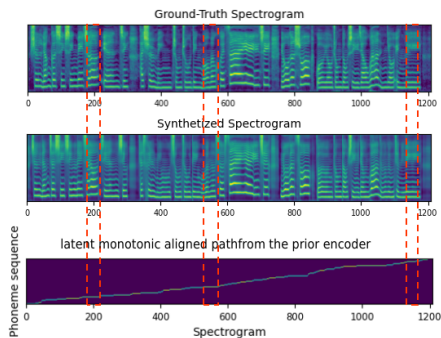


Figure: Visualization of GT, synthesized singing voice spectrogram, and latent optimal path from the prior encoder. *The GT and generated spectrogram are almost identical, and the generated spectrogram has a similar temporal structure to the inferred latent optimal path.*

Verify Behaviour of Latent Optimal Path in BDP-VAE

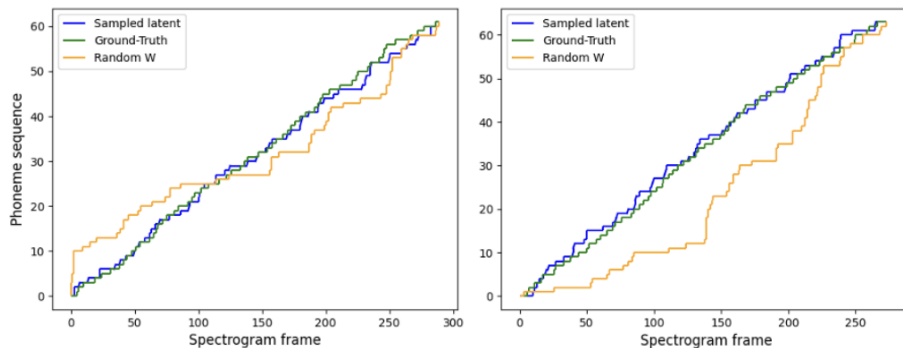


Figure: Visualization of GT alignment between phoneme tokens and spectrogram frames, latent optimal paths from the encoder, optimal paths from random latent space for two audio clips. *BDP-VAE achieves closer alignments with GT, indicating its effectiveness in finding latent optimal paths.*

Latent Optimal Paths by Gumbel Propagation for Variational Bayesian Dynamic Programming

Xinlei Niu¹, Christian Walder², Jing Zhang¹, Charles Patrick Martin¹

¹Australian National University

²Google DeepMind



Our Code



Our Paper