



Mitigating Modality Collapse in Multimodal VAEs via Impartial Optimization



Adrián Javaloy



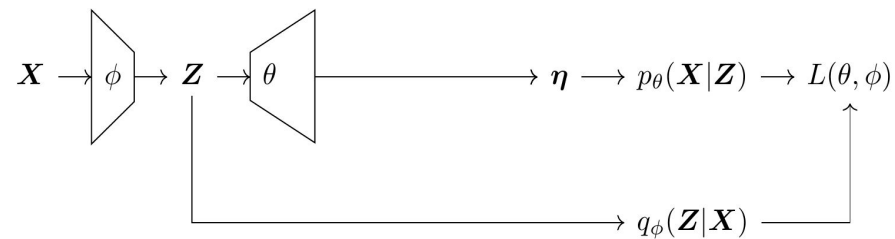
Maryam Meghdadi



Isabel Valera

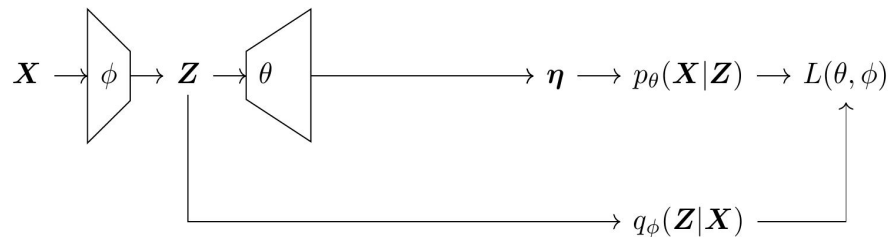
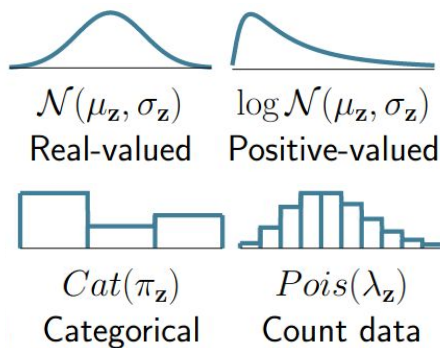
Problem motivation

Problem motivation



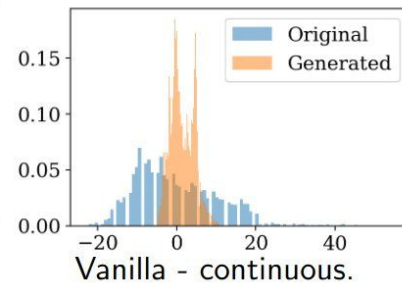
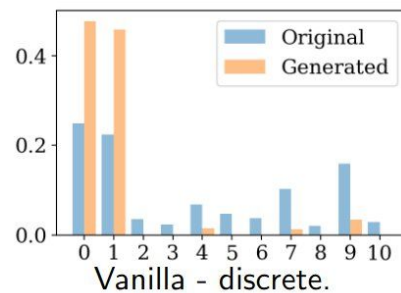
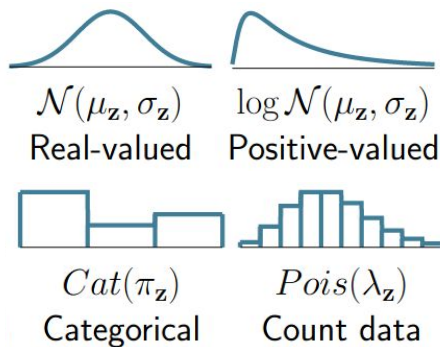
Problem motivation

Each modality is of a different type:



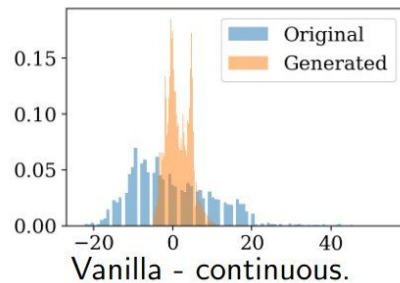
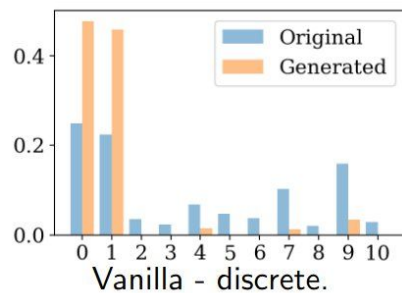
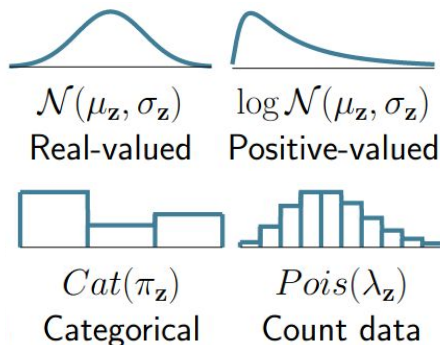
Problem motivation

Each modality is of a different type:



Problem motivation

Each modality is of a different type:

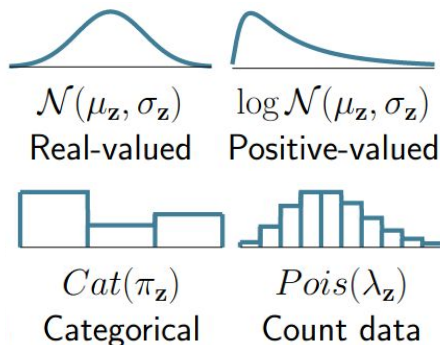


Likelihood impartiality

We aim for a learning process that does not prioritize learning any of the different likelihood modalities.

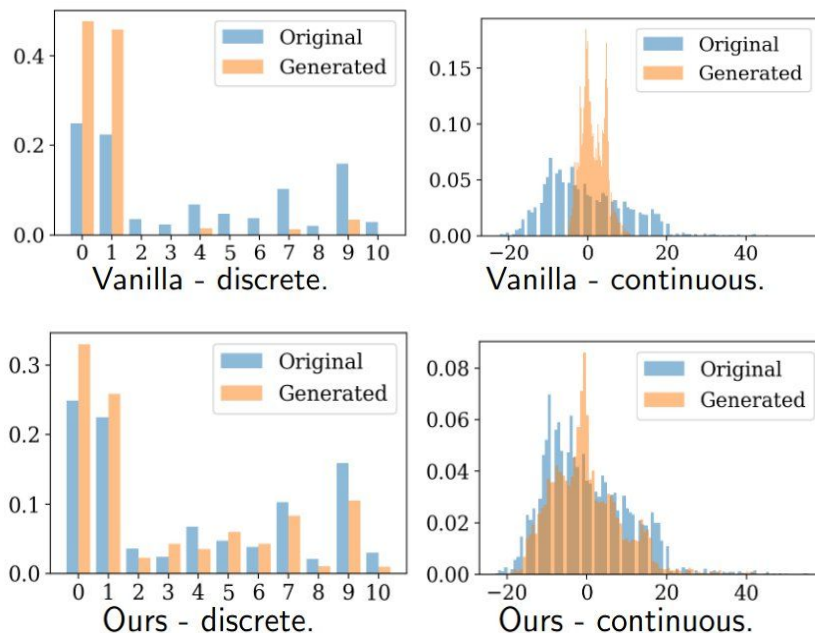
Problem motivation

Each modality is of a different type:



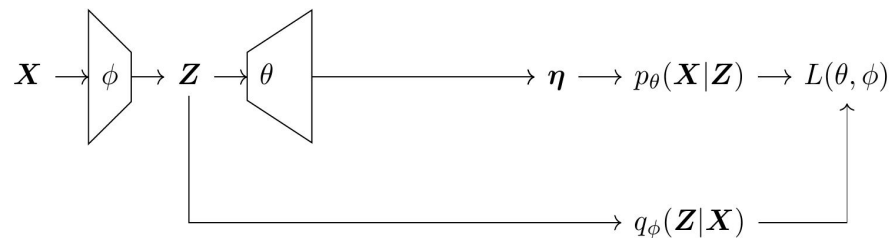
Likelihood impartiality

We aim for a learning process that does not prioritize learning any of the different likelihood modalities.

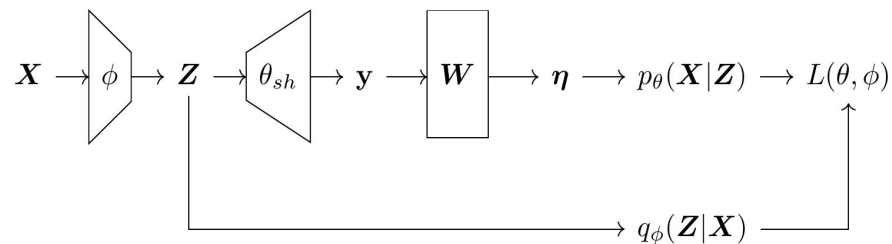


Analyzing modality collapse

Analyzing modality collapse

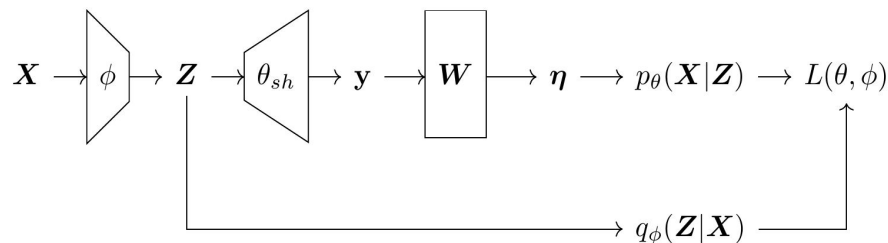
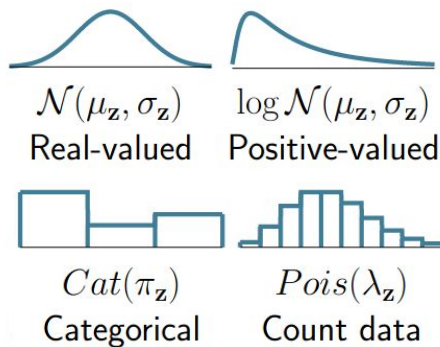


Analyzing modality collapse



Analyzing modality collapse

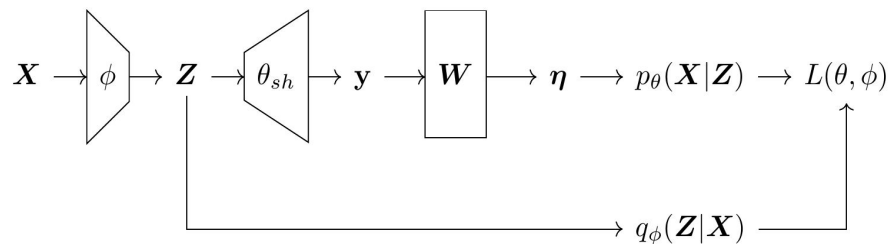
Each modality is of a different type:



Analyzing modality collapse

Each modality is of a different type:

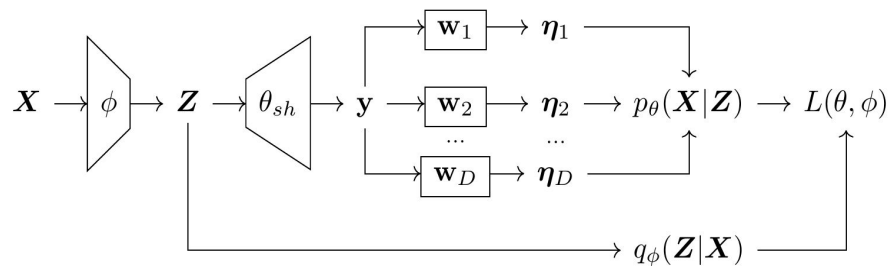
$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Analyzing modality collapse

Each modality is of a different type:

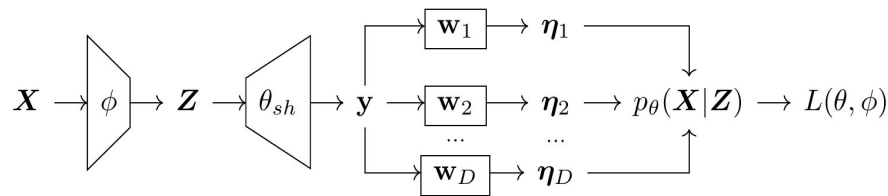
$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Analyzing modality collapse

Each modality is of a different type:

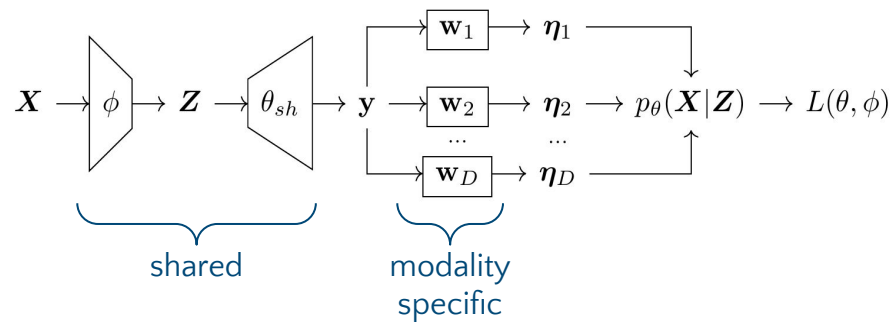
$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Analyzing modality collapse

Each modality is of a different type:

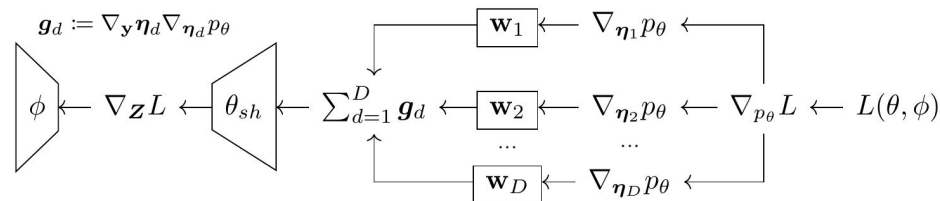
$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Analyzing modality collapse

Each modality is of a different type:

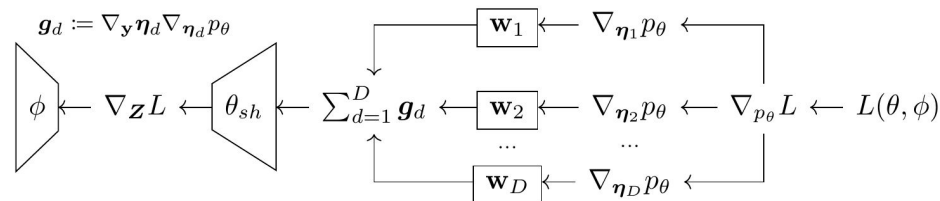
$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Analyzing modality collapse

Each modality is of a different type:

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$

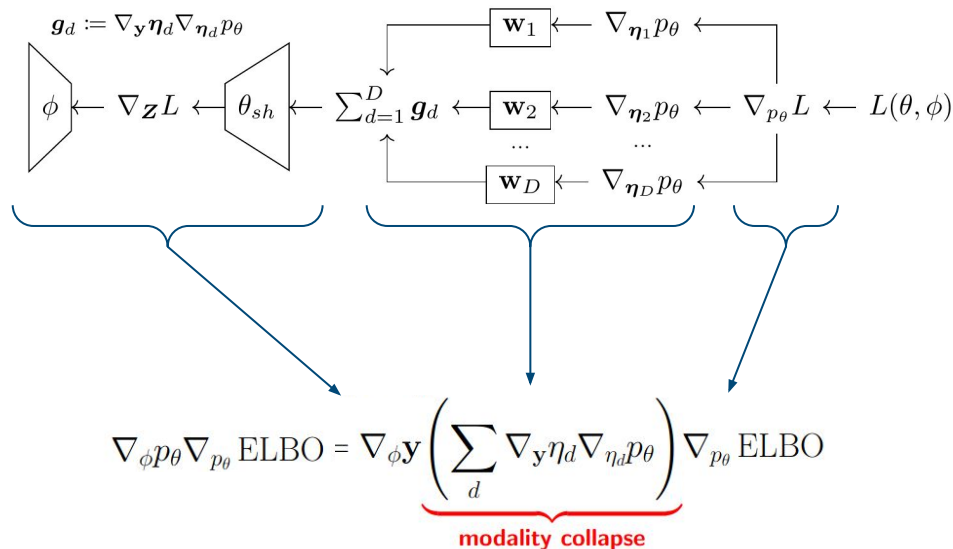


$$\nabla_{\phi} p_{\theta} \nabla_{p_{\theta}} \text{ELBO}$$

Analyzing modality collapse

Each modality is of a different type:

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$

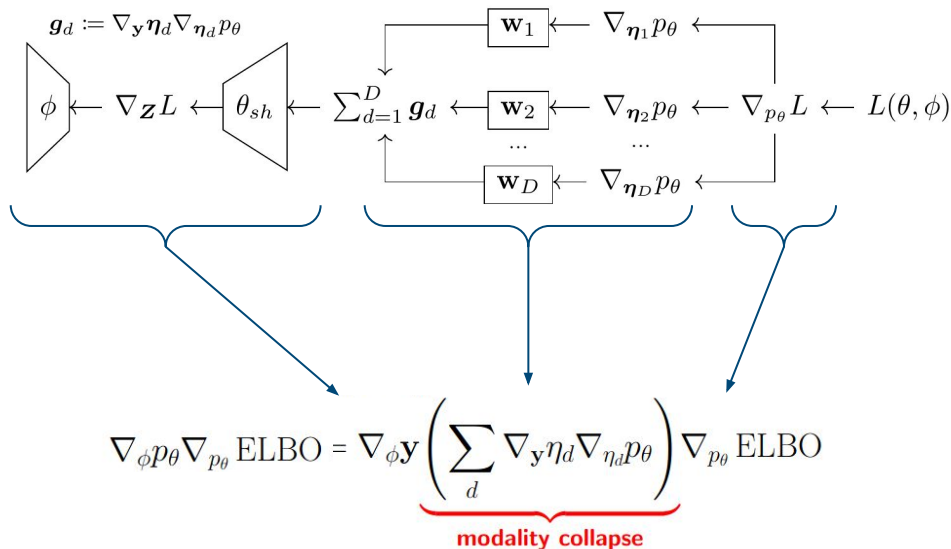


Analyzing modality collapse

Each modality is of a different type:

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$

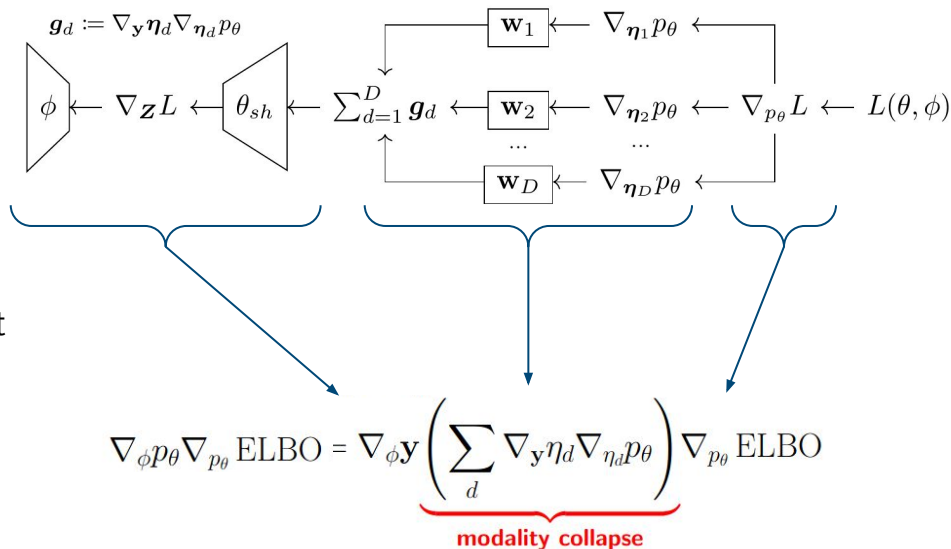
Modality collapse



Analyzing modality collapse

Each modality is of a different type:

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$

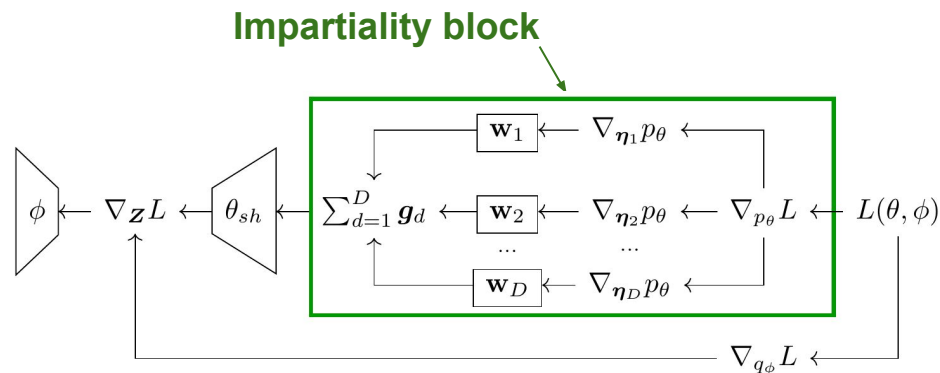


Modality collapse ← Gradient conflict

Analyzing modality collapse

Each modality is of a different type:

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Modality collapse ← Gradient conflict

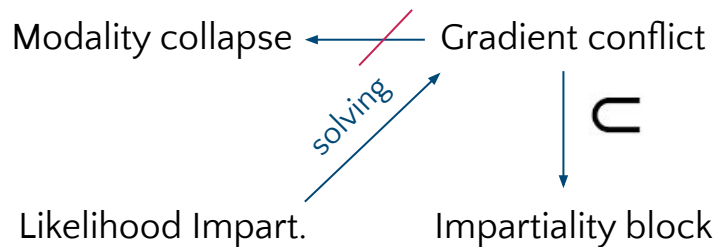
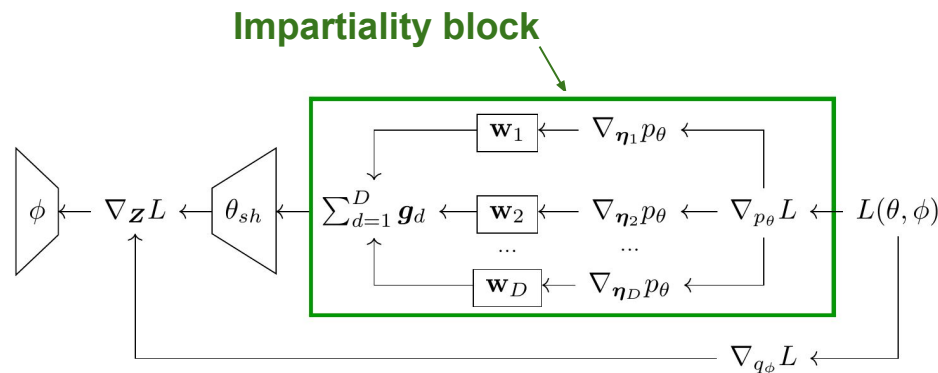
⊂

Impartiality block

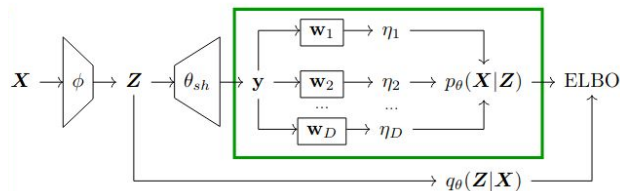
Analyzing modality collapse

Each modality is of a different type:

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$



Multimodal VAEs



Want: Model all modalities equally well.

Problem: Modality collapse.

Shared params: ϕ and θ_{sh} .

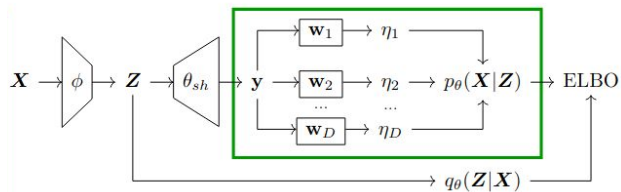
Excl. params: w_1, w_2, \dots, w_K .

Updating ϕ :

$$\begin{aligned} \nabla_{\phi} p_{\theta} \nabla_{p_{\theta}} \text{ELBO} &= \\ &= \nabla_{\phi} y \underbrace{\left(\sum_d \nabla_y \eta_d \nabla_{\eta_d} p_{\theta} \right)}_{\text{modality collapse}} \nabla_{p_{\theta}} \text{ELBO} \end{aligned}$$

A familiar face

Multimodal VAEs



Want: Model all modalities equally well.

Problem: Modality collapse.

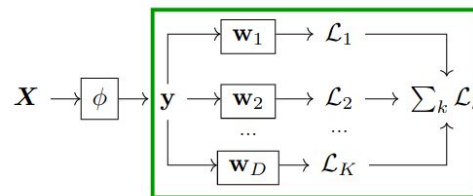
Shared params: ϕ and θ_{sh} .

Excl. params: w_1, w_2, \dots, w_K .

Updating ϕ :

$$\begin{aligned} \nabla_{\phi} p_{\theta} \nabla_{p_{\theta}} \text{ELBO} &= \\ &= \nabla_{\phi} \underbrace{\mathbf{y} \left(\sum_d \nabla_{\mathbf{y}} \eta_d \nabla_{\eta_d} p_{\theta} \right)}_{\text{modality collapse}} \nabla_{p_{\theta}} \text{ELBO} \end{aligned}$$

Multitask Learning



Want: Learn all tasks equally well.

Problem: Negative transfer.

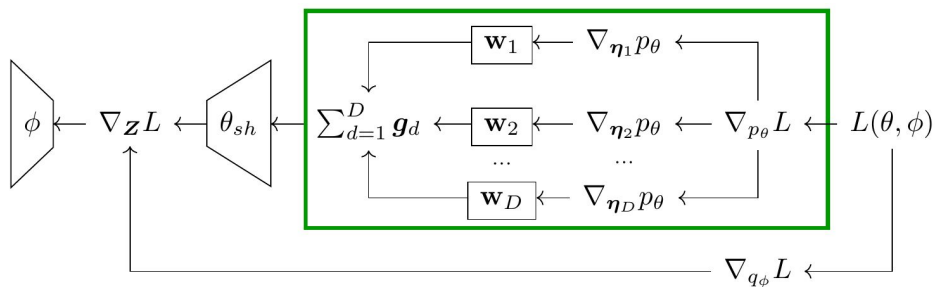
Shared params: ϕ .

Excl. params: w_1, w_2, \dots, w_K .

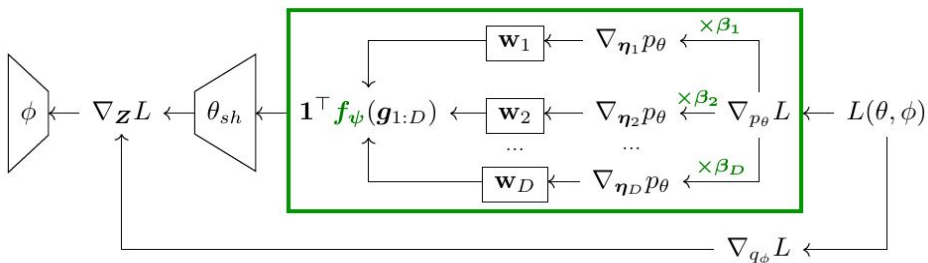
Updating ϕ :

$$\nabla_{\phi} \mathcal{L} = \nabla_{\phi} \mathbf{y} \underbrace{\left(\sum_k \nabla_{\mathbf{y}} \mathcal{L}_k \right)}_{\text{negative transfer}}$$

How to achieve impartiality



How to achieve impartiality

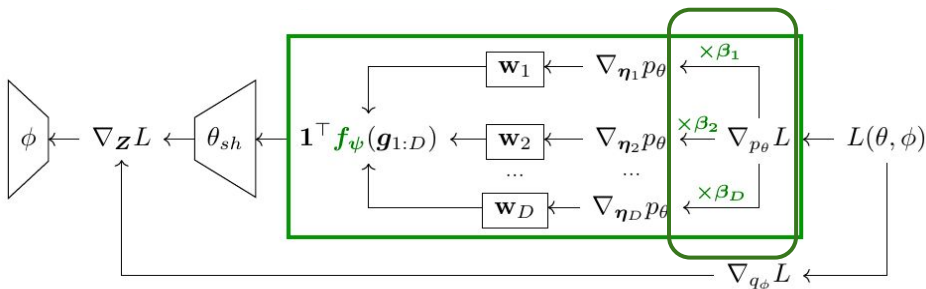


Two-step solution:

Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
- 2: **for** $d = 1$ **to** D **do**
- 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d} p_\theta \nabla_{p_\theta} L$
- 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
- 5: $\mathbf{g}_d \leftarrow \nabla_y \eta_d \cdot \mathbf{h}_d$
- 6: **end for**
- 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_\psi(\mathbf{g}_{1:D})$
- 8: **return** $\sum_d \tilde{\mathbf{g}}_d$

How to achieve impartiality



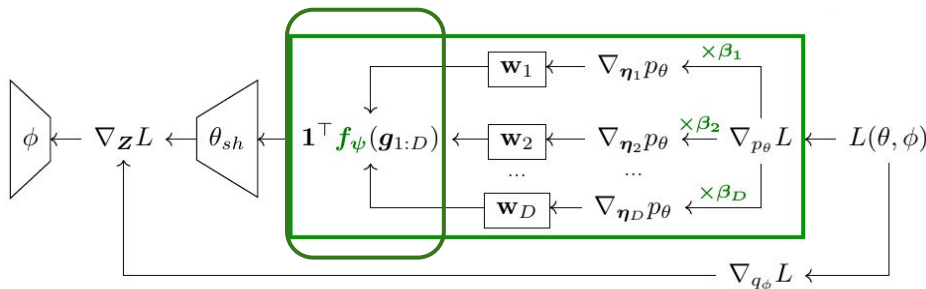
Two-step solution:

Local step, β : re-scale gradients to make them comparable.

Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
- 2: **for** $d = 1$ **to** D **do**
- 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d} p_\theta \nabla_{p_\theta} L$
- 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
- 5: $\mathbf{g}_d \leftarrow \nabla_{\mathbf{y}} \eta_d \cdot \mathbf{h}_d$
- 6: **end for**
- 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_\psi(\mathbf{g}_{1:D})$
- 8: **return** $\sum_d \tilde{\mathbf{g}}_d$

How to achieve impartiality



Algorithm 1 Backward pass within the impartiality block.

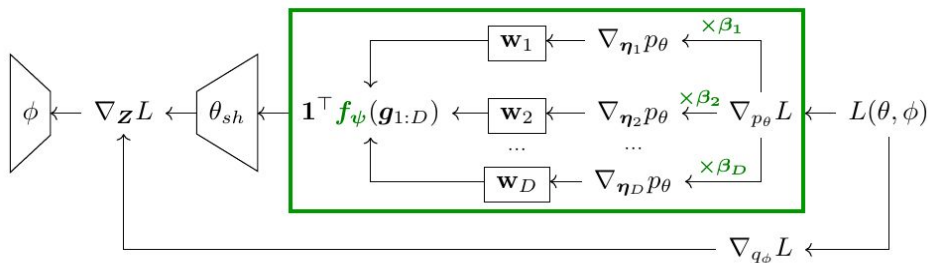
- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
- 2: **for** $d = 1$ **to** D **do**
- 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d} p_\theta \nabla_{p_\theta} L$
- 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
- 5: $\mathbf{g}_d \leftarrow \nabla_y \eta_d \cdot \mathbf{h}_d$
- 6: **end for**
- 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_\psi(\mathbf{g}_{1:D})$
- 8: **return** $\sum_d \mathbf{g}_d$

Two-step solution:

Local step, β : re-scale gradients to make them comparable.

Global step, \mathbf{f} : apply an MTL algorithm to alleviate gradient conflict.

How to achieve impartiality

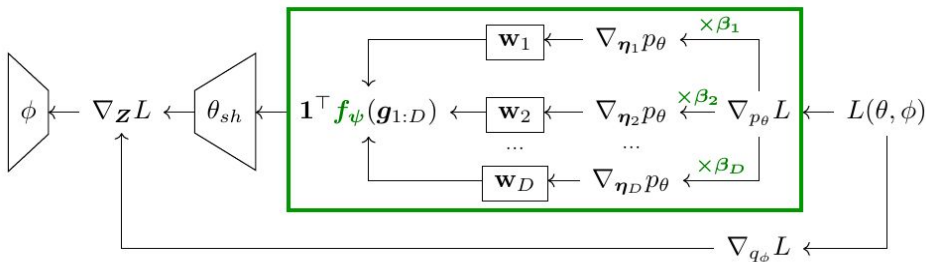


The impartiality block is local!

Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
- 2: **for** $d = 1$ **to** D **do**
- 3: $h_d \leftarrow \beta_d \nabla_{\eta_d} p_\theta \nabla_{p_\theta} L$
- 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot h_d$
- 5: $g_d \leftarrow \nabla_y \eta_d \cdot h_d$
- 6: **end for**
- 7: $\tilde{g}_{1:D} \leftarrow f_\psi(g_{1:D})$
- 8: **return** $\sum_d \tilde{g}_d$

How to achieve impartiality



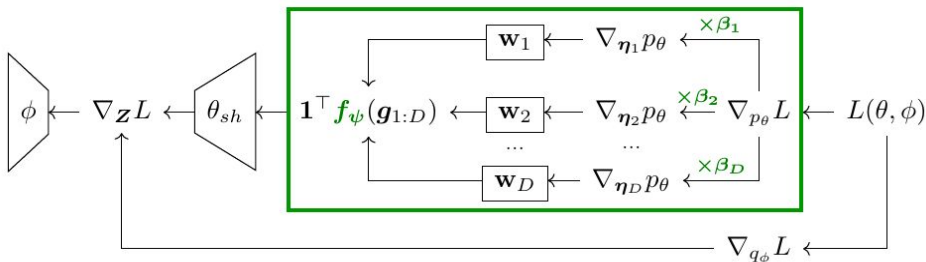
The impartiality block is local!

It can appear anywhere.

Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
 - 2: **for** $d = 1$ **to** D **do**
 - 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d} p_\theta \nabla_{p_\theta} L$
 - 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
 - 5: $\mathbf{g}_d \leftarrow \nabla_{\mathbf{y}} \eta_d \cdot \mathbf{h}_d$
 - 6: **end for**
 - 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_\psi(\mathbf{g}_{1:D})$
 - 8: **return** $\sum_d \tilde{\mathbf{g}}_d$
-

How to achieve impartiality



The impartiality block is local!

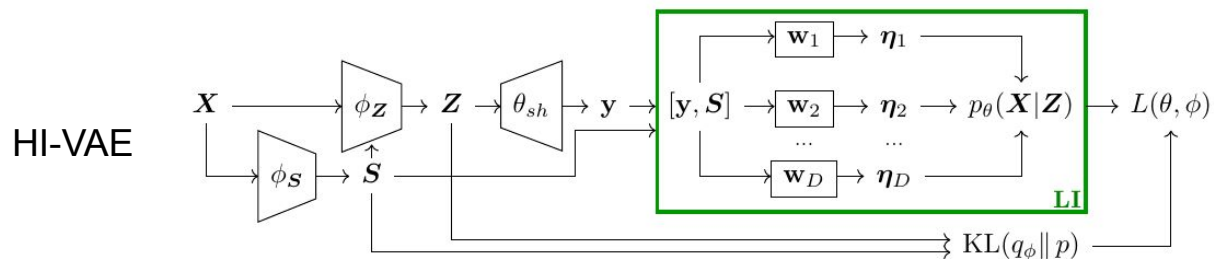
It can appear anywhere.

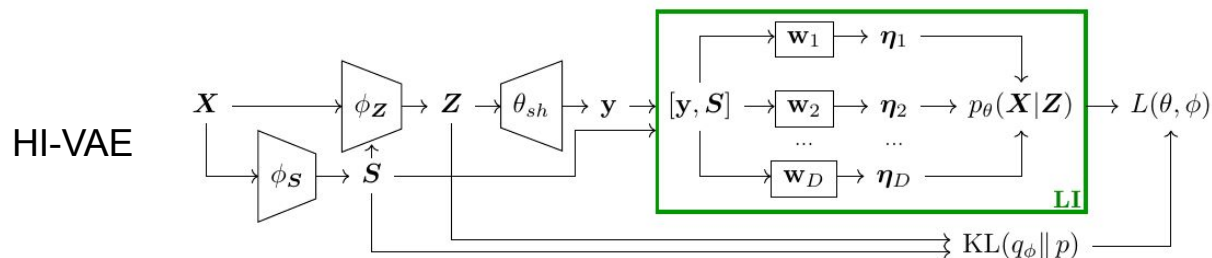
It can appear more than once.

Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_{\theta}} L$.
 - 2: **for** $d = 1$ **to** D **do**
 - 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d} p_{\theta} \nabla_{p_{\theta}} L$
 - 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
 - 5: $\mathbf{g}_d \leftarrow \nabla_{\mathbf{y}} \eta_d \cdot \mathbf{h}_d$
 - 6: **end for**
 - 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_{\psi}(\mathbf{g}_{1:D})$
 - 8: **return** $\sum_d \tilde{\mathbf{g}}_d$
-

VAE extensions

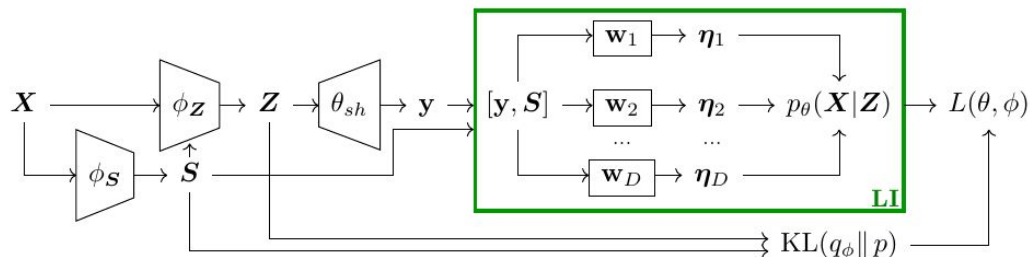




Two impartiality blocks:
One for y
One for S

VAE extensions

HI-VAE

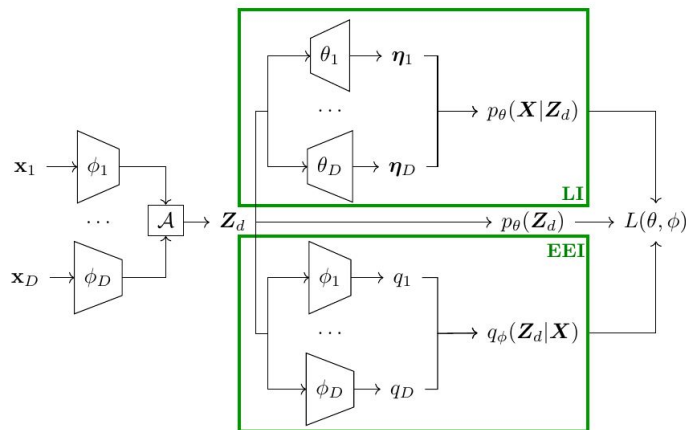


Two impartiality blocks:
One for y
One for S

MVAE

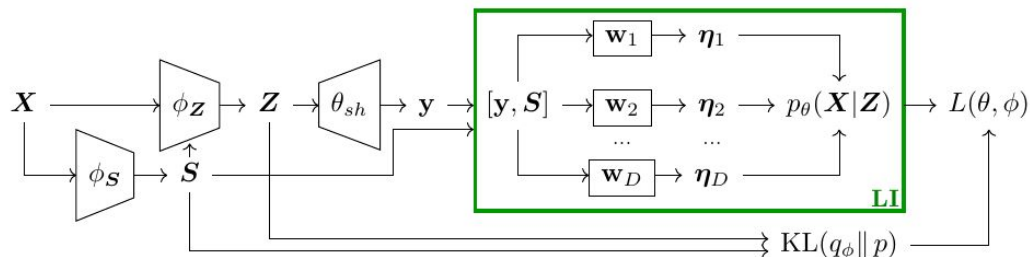
MMVAE

MoPoE



VAE extensions

HI-VAE

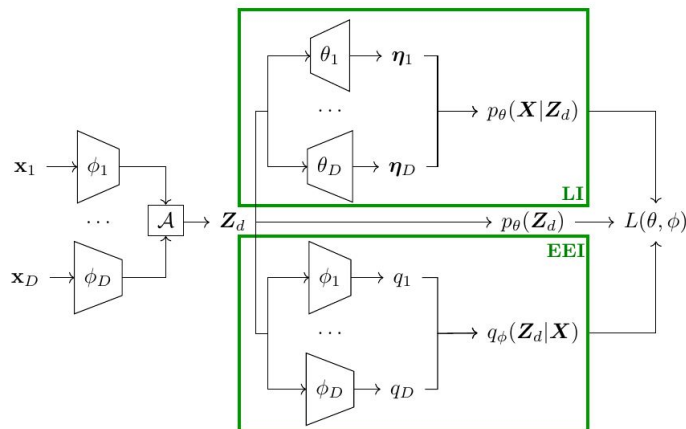


Two impartiality blocks:
One for y
One for S

MVAE

MMVAE

MoPoE



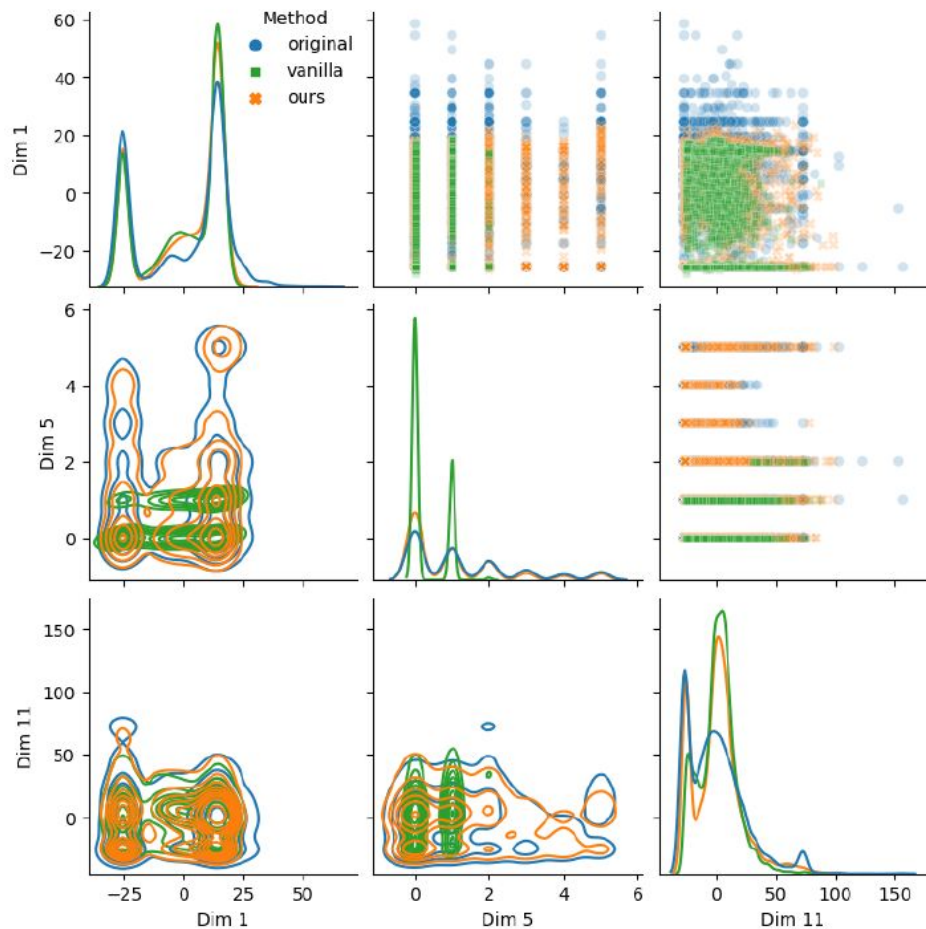
Additional goals:
Encoder Expert-Impartiality
Decoder Expert-Impartiality

As many impartiality blocks as modalities.

Results – heterogeneous data

Table 1. Test reconstruction errors (median over five seeds) for different datasets and VAE models. Statistically different values according to a corrected paired t-test ($\alpha = 0.1$) are shown in bold. Models trained with our approach outperforms the baseline in most cases.

			Heterogeneous									Homogeneous		
			<i>Adult</i>	<i>Credit</i>	<i>Wine</i>	<i>Diam.</i>	<i>Bank</i>	<i>IMDB</i>	<i>HI</i>	<i>rwm5yr</i>	<i>labour</i>	<i>El Nino</i>	<i>Magic</i>	<i>BooNE</i>
Standard VAE	ELBO	vanilla	0.213	0.128	0.086	0.187	0.203	0.082	0.170	0.105	0.109	0.109	0.064	0.042
		ours	0.104	0.041	0.071	0.139	0.043	0.032	0.041	0.026	0.063	0.068	0.058	0.039
	IWAE	vanilla	0.226	0.134	0.075	0.185	0.199	0.090	0.155	0.094	0.098	0.086	0.053	0.037
		ours	0.129	0.051	0.066	0.125	0.076	0.035	0.042	0.032	0.066	0.061	0.048	0.035
	DReG	vanilla	0.234	0.132	0.077	0.176	0.191	0.088	0.153	0.094	0.096	0.085	0.050	0.037
		ours	0.168	0.075	0.065	0.139	0.103	0.055	0.042	0.026	0.076	0.069	0.046	0.036
HI-VAE	vanilla	ours	0.127	0.107	0.126	0.114	0.141	0.079	0.105	0.044	0.100	0.098	0.062	0.039
		ours	0.081	0.060	0.117	0.011	0.095	0.049	0.109	0.024	0.069	0.015	0.033	0.038



Results – multimodal data

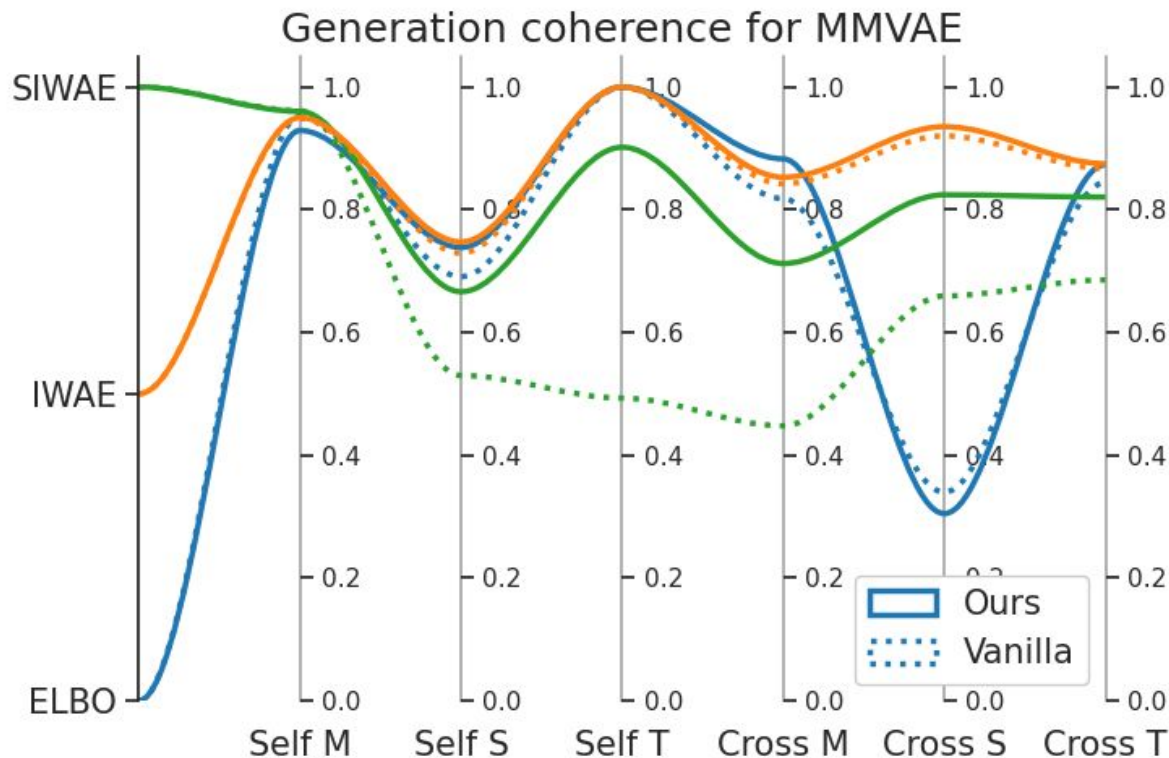
Table 5. Self and cross latent classification accuracy (%) for different models and losses on MNIST-SVHN-Text.

		ELBO	IWAE	SIWAE
		Self latent classification		
MVAE	vanilla	69.68	69.14	68.58
	ours	69.95	69.06	69.75
MMVAE	vanilla	71.81	87.55	71.30
	ours	87.83	90.78	85.55
MoPoE	vanilla	89.85	87.23	67.58
	ours	91.47	90.74	69.26
		Cross latent classification		
MVAE	vanilla	33.60	39.15	38.36
	ours	35.25	49.73	46.23
MMVAE	vanilla	44.25	76.81	40.60
	ours	71.42	84.80	60.50
MoPoE	vanilla	66.14	83.71	40.36
	ours	84.52	90.48	53.24

Table 9. Reconstruction coherence ($A = \{M, S, T\}$) for each modality, model, and dataset.

		ELBO			IWAE			SIWAE		
	\mathbf{x}_d	M	S	T	M	S	T	M	S	T
MVAE	vanilla	97.53	88.26	99.30	97.27	87.19	98.76	97.37	87.47	98.83
	ours	97.85	89.65	99.64	98.28	89.01	99.93	97.42	87.63	99.20
MMVAE	vanilla	86.01	45.59	89.17	85.25	84.03	88.66	58.95	61.27	63.27
	ours	89.42	45.83	91.54	87.55	86.87	90.93	74.85	73.89	81.09
MoPoE	vanilla	95.72	85.86	98.01	95.82	87.55	97.93	75.10	67.16	76.61
	ours	96.50	93.60	99.14	97.29	92.93	99.00	96.91	89.01	99.28

Results – multimodal data



Make sure to follow us!



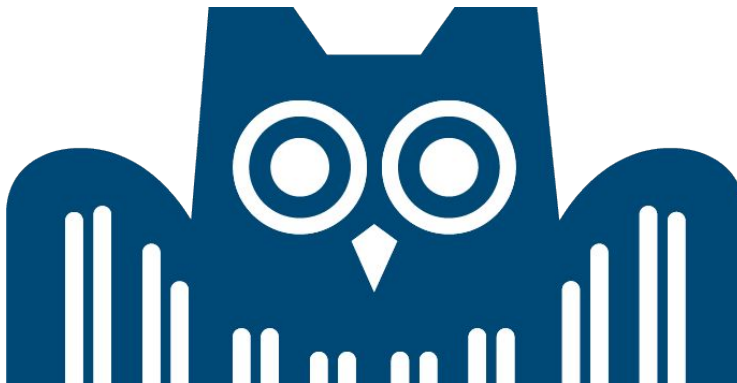
@javaloyML



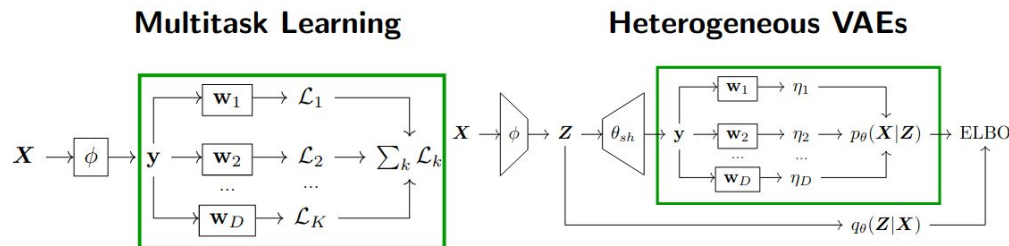
@maryam_meghdadi



@IValeraM



Gradient conflict in multimodal VAEs



Want: Learn all tasks equally well. **Want:** Model all features equally well.

Problem: Negative transfer.

Problem: Feature overlooking.

Shared params: ϕ .

Shared params: ϕ and θ_{sh} .

Excl. params: w_1, w_2, \dots, w_K .

Excl. params: w_1, w_2, \dots, w_K .

Updating ϕ :

Updating ϕ :

$$\nabla_{\phi} \mathcal{L} = \nabla_{\phi} \mathbf{y} \left(\underbrace{\sum_k \nabla_{\mathbf{y}} \mathcal{L}_k}_{\text{negative transfer}} \right)$$

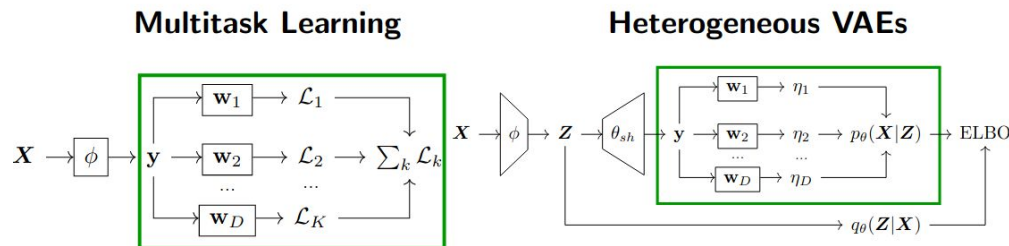
$$\begin{aligned} \nabla_{\phi} p_{\theta} \nabla_{p_{\theta}} \text{ELBO} &= \\ &= \nabla_{\phi} \mathbf{y} \left(\underbrace{\sum_d \nabla_{\mathbf{y}} \eta_d \nabla_{\eta_d} p_{\theta}}_{\text{feature overlooking}} \right) \nabla_{p_{\theta}} \text{ELBO} \end{aligned}$$

Table 3. Reconstruction coherence ($A = \{M, S, T\}$) for each modality and model, trained using SIWAE.

	x_d	M	S	T
MVAE	vanilla	97.37	87.47	98.83
	ours	97.42	87.63	99.20
MMVAE	vanilla	58.95	61.27	63.27
	ours	74.16	68.93	78.17
MoPoE	vanilla	75.10	67.16	76.61
	ours	96.91	89.01	99.28

			<i>Adult</i>	<i>Credit</i>	<i>Wine</i>
Standard VAE	ELBO	vanilla	0.213	0.128	0.086
		ours	0.104	0.041	0.071
	IWAE	vanilla	0.226	0.134	0.075
		ours	0.129	0.051	0.066
	DReG	vanilla	0.234	0.132	0.077
		ours	0.168	0.075	0.065
HI-VAE	vanilla	0.127	0.107	0.126	
	ours	0.081	0.060	0.117	

Gradient conflict in multimodal VAEs



Want: Learn all tasks equally well. **Want:** Model all features equally well.

Problem: Negative transfer.

Problem: Feature overlooking.

Shared params: ϕ .

Shared params: ϕ and θ_{sh} .

Excl. params: w_1, w_2, \dots, w_K .

Excl. params: w_1, w_2, \dots, w_K .

Updating ϕ :

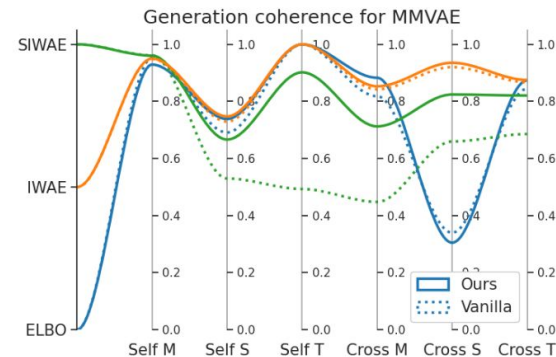
Updating ϕ :

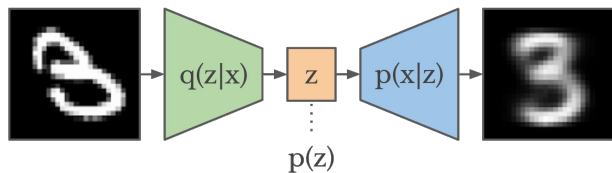
$$\nabla_{\phi} \mathcal{L} = \nabla_{\phi} \mathbf{y} \left(\underbrace{\sum_k \nabla_{\mathbf{y}} \mathcal{L}_k}_{\text{negative transfer}} \right)$$

$$\begin{aligned} \nabla_{\phi} p_{\theta} \nabla_{p_{\theta}} \text{ELBO} &= \\ &= \nabla_{\phi} \mathbf{y} \left(\underbrace{\sum_d \nabla_{\mathbf{y}} \eta_d \nabla_{\eta_d} p_{\theta}}_{\text{feature overlooking}} \right) \nabla_{p_{\theta}} \text{ELBO} \end{aligned}$$

Table 3. Reconstruction coherence ($A = \{M, S, T\}$) for each modality and model, trained using SIWAE.

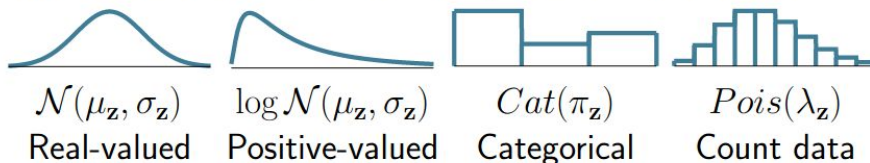
	x_d	M	S	T
MVAE	vanilla	97.37	87.47	98.83
	ours	97.42	87.63	99.20
MMVAE	vanilla	58.95	61.27	63.27
	ours	74.16	68.93	78.17
MoPoE	vanilla	75.10	67.16	76.61
	ours	96.91	89.01	99.28





$$\underset{\theta, \phi}{\text{maximize}} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))$$

- **Heterogeneous data.** Each feature \mathbf{x}_d is of a different type:



Thus, a heterogeneous likelihood is of the form

$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z}).$$

Likelihood impartiality



All dimensions are equally important. We want a learning process impartial to the likelihood of each of the dimensions.

We noticed

Each modality is of a different type:

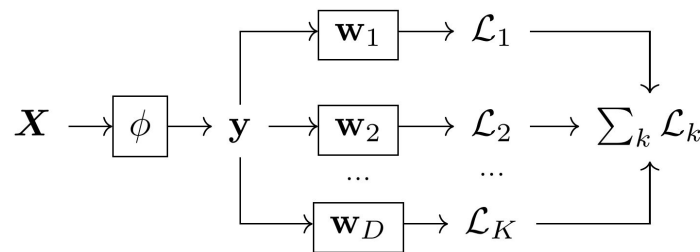
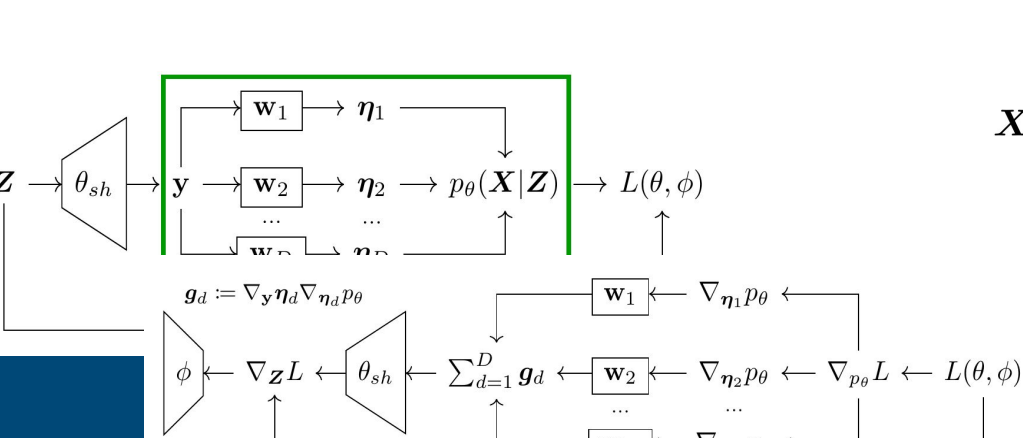
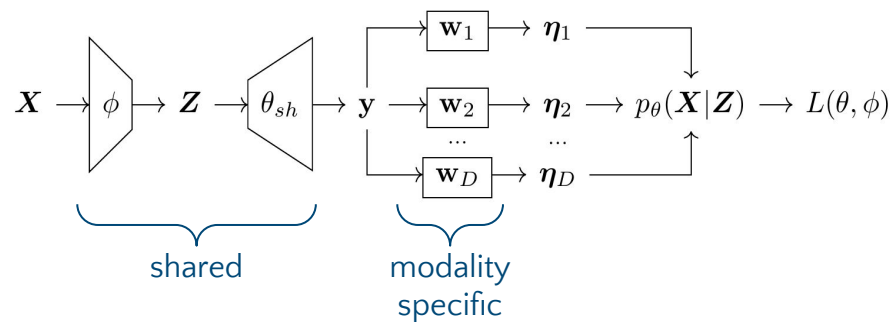
$$p_{\theta}(\mathbf{X}|\mathbf{z}) = \prod_{d=1}^D p_d(\mathbf{x}_d|\mathbf{z})$$

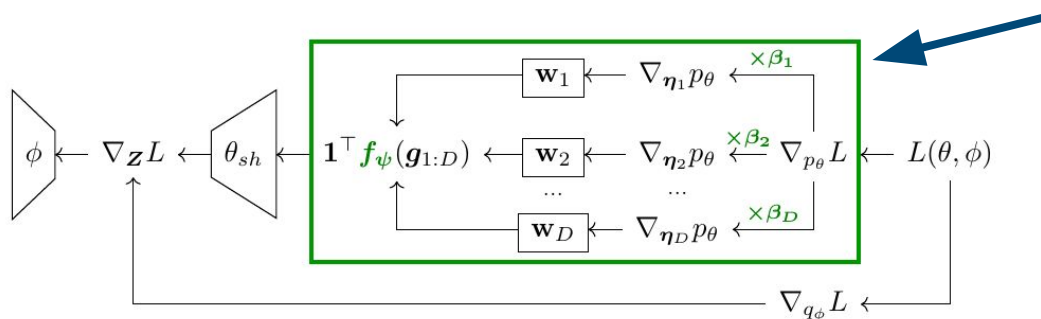
Real-valued Positive-valued

Cat($\pi_{\mathbf{z}}$) Pois($\lambda_{\mathbf{z}}$)

Categorical Count data





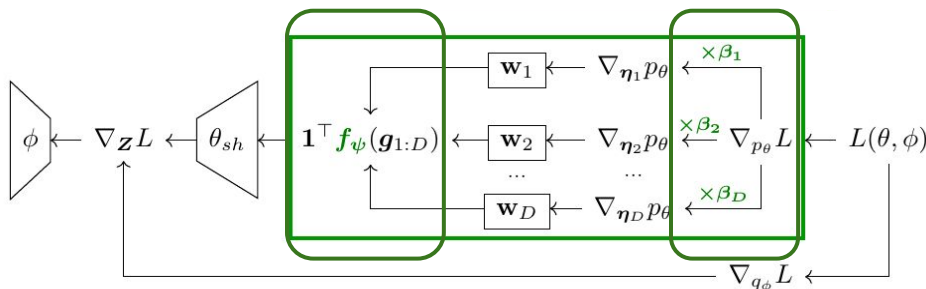
Impartiality block

- Input: shared
- Split-and-merge structure
- Suffers gradient conflict
- Two-step solution:
 - Local step, β : re-scale grads to make them comparable.
 - Global step, f : apply an MTL algorithm to alleviate gradient conflict.
- *Local* character
 - Anywhere
 - More than once

Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
- 2: **for** $d = 1$ **to** D **do**
- 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d} p_\theta \nabla_{p_\theta} L$
- 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
- 5: $\mathbf{g}_d \leftarrow \nabla_{\mathbf{y}} \eta_d \cdot \mathbf{h}_d$
- 6: **end for**
- 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_\psi(\mathbf{g}_{1:D})$
- 8: **return** $\sum_d \tilde{\mathbf{g}}_d$

How to achieve impartiality



Algorithm 1 Backward pass within the impartiality block.

- 1: **Input:** Output gradient, $\nabla_{p_\theta} L$.
- 2: **for** $d = 1$ **to** D **do**
- 3: $\mathbf{h}_d \leftarrow \beta_d \nabla_{\eta_d p_\theta} \nabla_{p_\theta} L$
- 4: $\nabla_{\omega_d} L \leftarrow \nabla_{\omega_d} \eta_d \cdot \mathbf{h}_d$
- 5: $\mathbf{g}_d \leftarrow \nabla_y \eta_d \cdot \mathbf{h}_d$
- 6: **end for**
- 7: $\tilde{\mathbf{g}}_{1:D} \leftarrow \mathbf{f}_\psi(\mathbf{g}_{1:D})$
- 8: **return** $\sum_d \mathbf{g}_d$

The impartiality block is local!

Local/slope preserving gradients to make the solution as close as possible

Doesn't depend on different modalities to alleviate gradient conflict.

Analyzing modality collapse

Each modality is of a different type:

