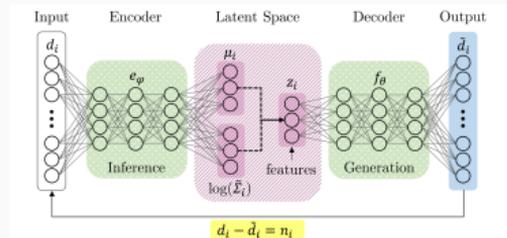


# Probabilistic Variational Autoencoders

## A BAYESIAN PERSPECTIVE



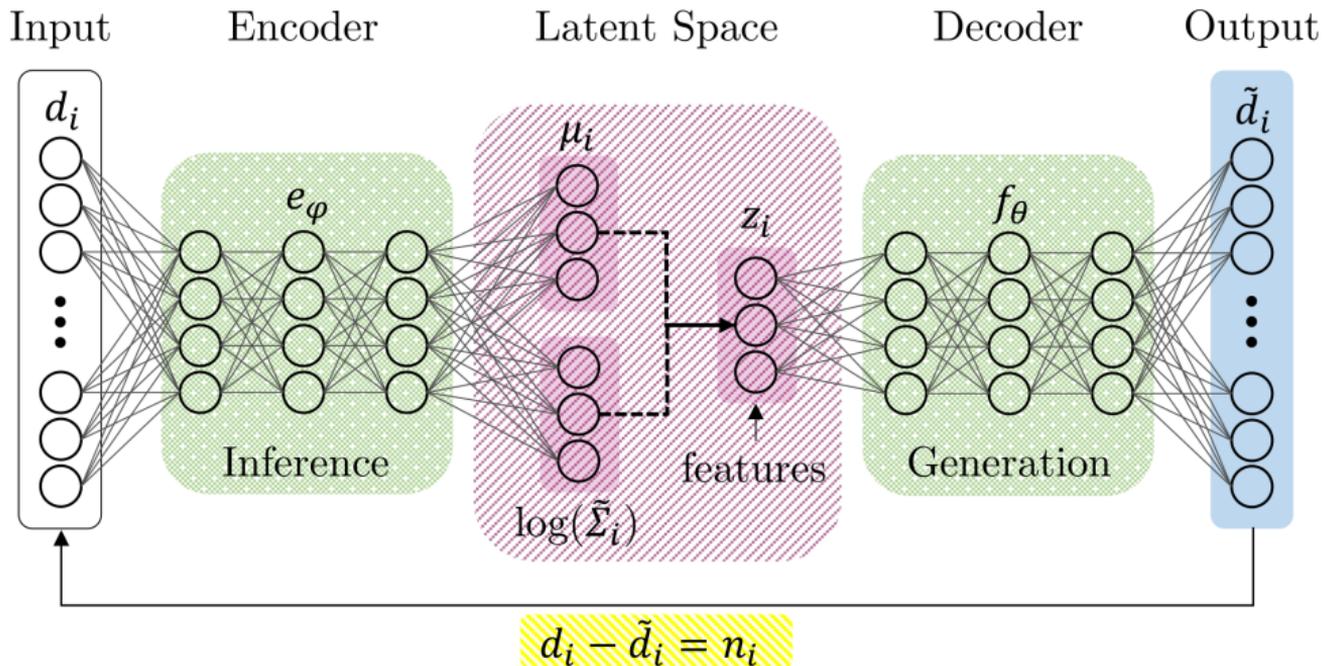
Philipp Frank

February 7, 2022

Berlin Machine Learning Group, Germany

[www.ph-frank.de](http://www.ph-frank.de)

# Variational Autoencoder (VAE)



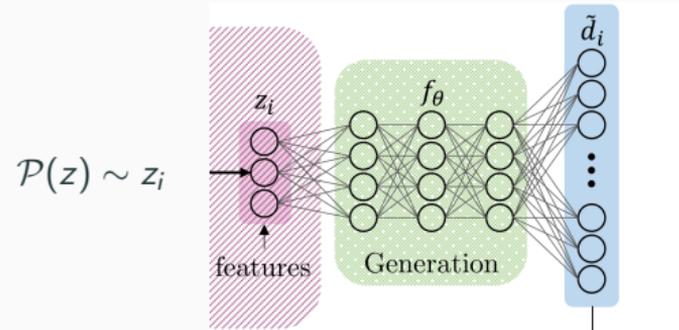
# Probabilistic generative model

---

# VAE - Generative model (Decoder)

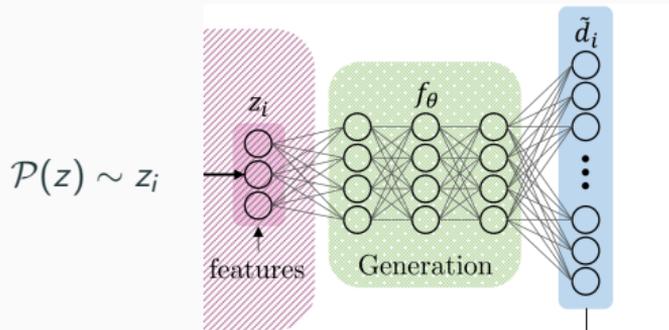
- † Latent features from common distribution

$$\mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$$



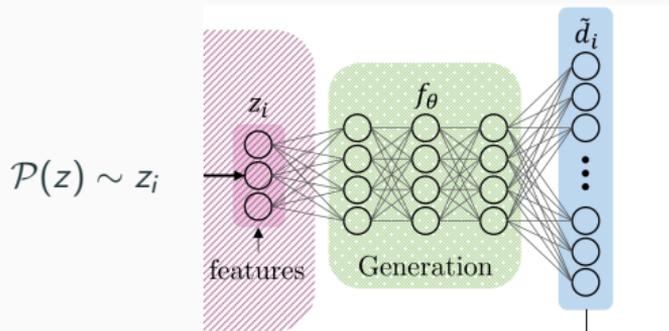
# VAE - Generative model (Decoder)

- † Latent features from common distribution  
 $\mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$
- † Generative model  $\tilde{d} = f_{\theta}(z)$



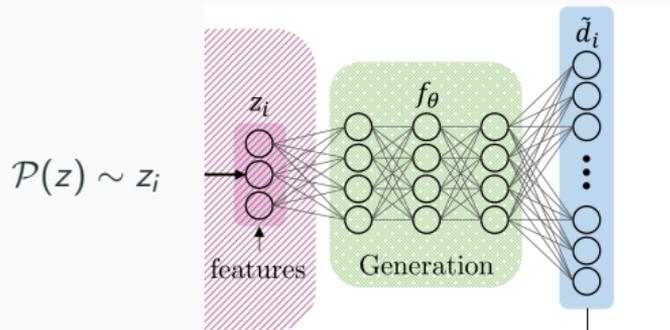
# VAE - Generative model (Decoder)

- † Latent features from common distribution  
 $\mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$
- † Generative model  $\tilde{d} = f_{\theta}(z)$
- † Probabilistic model  $d = \tilde{d} + n; \quad n \sim \mathcal{P}(n)$



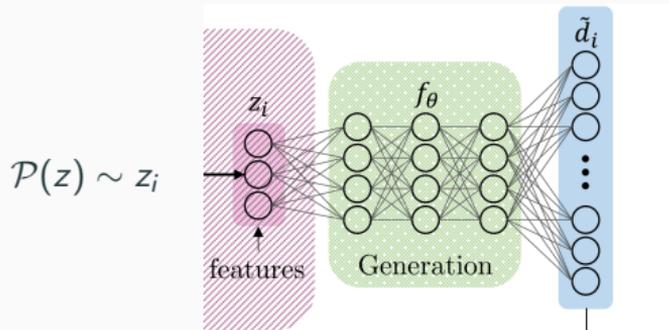
# VAE - Generative model (Decoder)

- † Latent features from common distribution  
 $\mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$
- † Generative model  $\tilde{d} = f_{\theta}(z)$
- † Probabilistic model  $d = \tilde{d} + n; \quad n \sim \mathcal{P}(n)$
- † E.g.  $\mathcal{P}(n) = \mathcal{N}(n|0, \mathbb{1})$



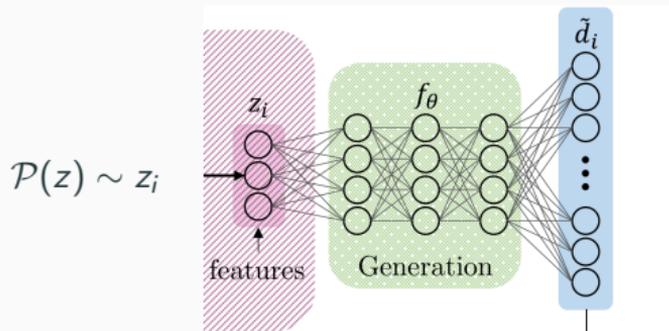
# VAE - Generative model (Decoder)

- † Latent features from common distribution  
 $\mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$
- † Generative model  $\tilde{d} = f_{\theta}(z)$
- † Probabilistic model  $d = \tilde{d} + n; \quad n \sim \mathcal{P}(n)$
- † E.g.  $\mathcal{P}(n) = \mathcal{N}(n|0, N)$



# VAE - Generative model (Decoder)

- † Latent features from common distribution  
 $\mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$
- † Generative model  $\tilde{d} = f_{\theta}(z)$
- † Probabilistic model  $d = \tilde{d} + n; \quad n \sim \mathcal{P}(n)$
- † E.g.  $\mathcal{P}(n) = \mathcal{N}(n|0, N)$
- † Adaptive covariance  $N \sim \mathcal{P}(N)$



# Inference

---

- + Likelihood  $\mathcal{P}(d_i|\theta, z_i) = \mathcal{N}(d_i|f_\theta(z_i), N)$
- + Prior  $\mathcal{P}(\theta, z_i) = \mathcal{P}(\theta) \mathcal{N}(z_i|0, \mathbf{1})$

- + Likelihood  $\mathcal{P}(d_i|\theta, z_i) = \mathcal{N}(d_i|f_\theta(z_i), N)$
- + Prior  $\mathcal{P}(\theta, z_i) = \mathcal{P}(\theta) \mathcal{N}(z_i|0, \mathbf{1})$
- + Dataset  $\mathcal{D} \equiv \{d_i\}_{i \in \{1, \dots, M\}}$ ; Latent variables  $\mathcal{Z} \equiv \{z_i\}_{i \in \{1, \dots, M\}}$

- + Likelihood  $\mathcal{P}(d_i|\theta, z_i) = \mathcal{N}(d_i|f_\theta(z_i), N)$
- + Prior  $\mathcal{P}(\theta, z_i) = \mathcal{P}(\theta) \mathcal{N}(z_i|0, \mathbb{1})$
- + Dataset  $\mathcal{D} \equiv \{d_i\}_{i \in \{1, \dots, M\}}$ ; Latent variables  $\mathcal{Z} \equiv \{z_i\}_{i \in \{1, \dots, M\}}$

## Product Rule of Probabilities aka **Bayes' theorem**

$$\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D}) = \frac{\left[ \prod_{i=1}^M \mathcal{P}(d_i|\theta, z_i) \mathcal{P}(z_i) \right] \mathcal{P}(\theta)}{\mathcal{P}(\mathcal{D})} .$$

- + Likelihood  $\mathcal{P}(d_i|\theta, z_i) = \mathcal{N}(d_i|f_\theta(z_i), N)$
- + Prior  $\mathcal{P}(\theta, z_i) = \mathcal{P}(\theta) \mathcal{N}(z_i|0, \mathbb{1})$
- + Dataset  $\mathcal{D} \equiv \{d_i\}_{i \in \{1, \dots, M\}}$ ; Latent variables  $\mathcal{Z} \equiv \{z_i\}_{i \in \{1, \dots, M\}}$

## Product Rule of Probabilities aka **Bayes' theorem**

$$\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D}) = \frac{\left[ \prod_{i=1}^M \mathcal{N}(d_i|f_\theta(z_i), N) \mathcal{N}(z_i|0, \mathbb{1}) \right] \mathcal{P}(\theta)}{\mathcal{P}(\mathcal{D})} .$$

- + Likelihood  $\mathcal{P}(d_i|\theta, z_i, N) = \mathcal{N}(d_i|f_\theta(z_i), N)$
- + Prior  $\mathcal{P}(\theta, z_i, N) = \mathcal{P}(\theta) \mathcal{N}(z_i|0, \mathbb{1}) \mathcal{P}(N)$
- + Dataset  $\mathcal{D} \equiv \{d_i\}_{i \in \{1, \dots, M\}}$ ; Latent variables  $\mathcal{Z} \equiv \{z_i\}_{i \in \{1, \dots, M\}}$

## Product Rule of Probabilities aka **Bayes' theorem**

$$\mathcal{P}(\theta, \mathcal{Z}, N|\mathcal{D}) = \frac{\left[ \prod_{i=1}^M \mathcal{N}(d_i|f_\theta(z_i), N) \mathcal{N}(z_i|0, \mathbb{1}) \right] \mathcal{P}(\theta) \mathcal{P}(N)}{\mathcal{P}(\mathcal{D})} .$$

# Variational Approximation

+ Posterior distribution  $\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})$  intractable

# Variational Approximation

- + Posterior distribution  $\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})$  intractable
- + Approximation with tractable distribution  $Q(\theta, \mathcal{Z}|\mathcal{D})$

- + Posterior distribution  $\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})$  intractable
- + Approximation with tractable distribution  $Q(\theta, \mathcal{Z}|\mathcal{D})$

## Kullback-Leibler Divergence

$$\mathcal{KL}[Q; \mathcal{P}] \equiv - \int \log \left( \frac{\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})}{Q(\theta, \mathcal{Z}|\mathcal{D})} \right) Q(\theta, \mathcal{Z}|\mathcal{D}) \, d\theta \, d\mathcal{Z} .$$

# Variational Approximation

- + Posterior distribution  $\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})$  intractable
- + Approximation with tractable distribution  $Q(\theta, \mathcal{Z}|\mathcal{D})$

## Kullback-Leibler Divergence

$$\mathcal{KL}[Q; \mathcal{P}] \equiv - \int \log \left( \frac{\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})}{Q(\theta, \mathcal{Z}|\mathcal{D})} \right) Q(\theta, \mathcal{Z}|\mathcal{D}) \, d\theta \, d\mathcal{Z} .$$

- + Parameterize approximation:  $Q(\theta, \mathcal{Z}|\mathcal{D}) = \delta(\theta - \bar{\theta}) \prod_{i=1}^M Q_{\phi}(z_i|d_i)$

# Variational Approximation

- + Posterior distribution  $\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})$  intractable
- + Approximation with tractable distribution  $Q(\theta, \mathcal{Z}|\mathcal{D})$

## Kullback-Leibler Divergence

$$\mathcal{KL}[Q; \mathcal{P}] \equiv - \int \log \left( \frac{\mathcal{P}(\theta, \mathcal{Z}|\mathcal{D})}{Q(\theta, \mathcal{Z}|\mathcal{D})} \right) Q(\theta, \mathcal{Z}|\mathcal{D}) \, d\theta \, d\mathcal{Z} .$$

- + Parameterize approximation:  $Q(\theta, \mathcal{Z}|\mathcal{D}) = \delta(\theta - \bar{\theta}) \prod_{i=1}^M Q_{\phi}(z_i|d_i)$
- + Optimize the  $\mathcal{KL}$  for  $\bar{\theta}$  and  $\phi$

## Inference Network (Encoder)

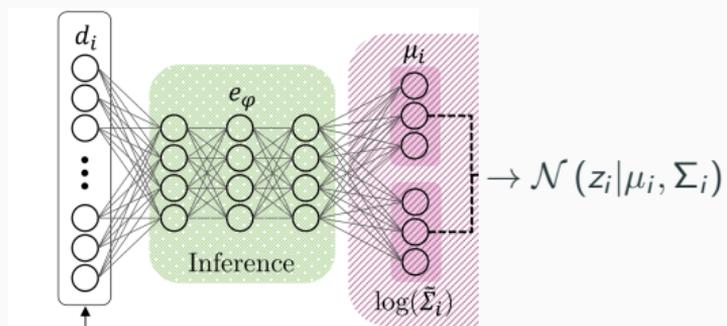
- † Parameterize as a normal distribution

$$Q_{\phi}(z_i|d_i) \equiv \mathcal{N}(z_i|\mu_i, \Sigma_i)$$

$$\rightarrow \mathcal{N}(z_i|\mu_i, \Sigma_i)$$

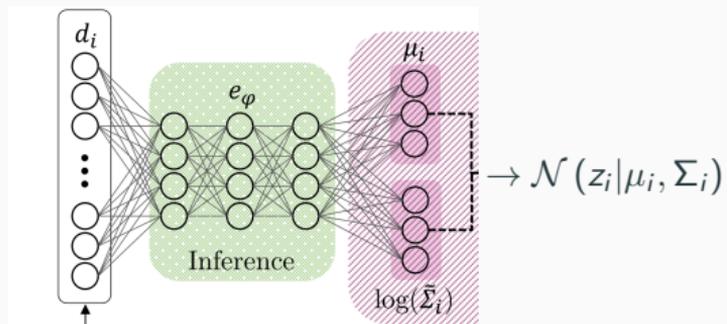
# Inference Network (Encoder)

- Parameterize as a normal distribution  
 $Q_\phi(z_i|d_i) \equiv \mathcal{N}(z_i|\mu_i, \Sigma_i)$
- And inference network  $\begin{pmatrix} \mu_i \\ \Sigma_i \end{pmatrix} = e_\phi(d_i)$



# Inference Network (Encoder)

- † Parameterize as a normal distribution  
 $Q_\phi(z_i|d_i) \equiv \mathcal{N}(z_i|\mu_i, \Sigma_i)$
- † And inference network  $\begin{pmatrix} \mu_i \\ \Sigma_i \end{pmatrix} = e_\phi(d_i)$
- †  $\Sigma$  often parameterized as a diagonal matrix



## **Extensions / Modifications**

---

- + Exchange / Modify likelihood distribution

- + Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...

- + Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...
- + Alter approximation distribution  $Q$

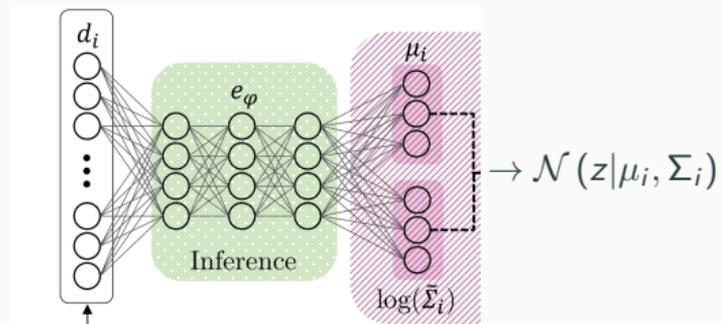
- + Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...
- + Alter approximation distribution  $Q$ 
  - Use full covariance matrix

- + Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...
- + Alter approximation distribution  $Q$ 
  - Use full covariance matrix
  - Estimate covariance from generative model using Fisher Information metric
  - ...

# Fisher metric as covariance

## Fisher Information metric

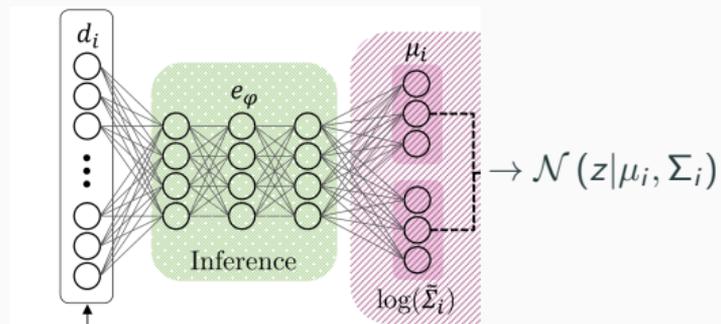
$$\mathcal{M}(\mu) \equiv - \left\langle \frac{\partial^2 \log(\mathcal{P})(d|z, \theta)}{\partial z \partial z} \right\rangle_{\mathcal{P}(d|z, \theta)} \Big|_{z=\mu} + \mathbf{1}$$



# Fisher metric as covariance

## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right)^{\dagger} N^{-1} \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right) \right]_{\mathbf{z}=\mu} + \mathbf{1}$$

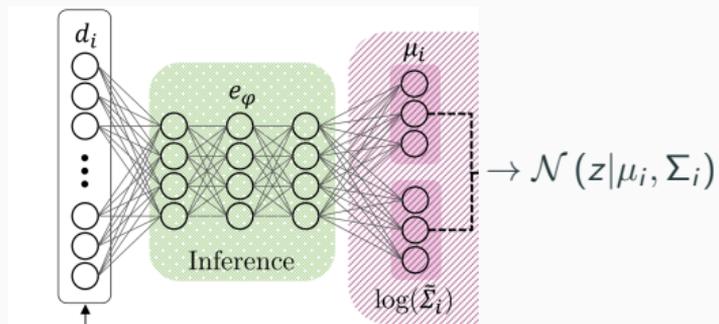


## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right)^{\dagger} N^{-1} \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right) \right]_{\mathbf{z}=\mu} + \mathbb{1}$$

- † Kramer-Rao bound

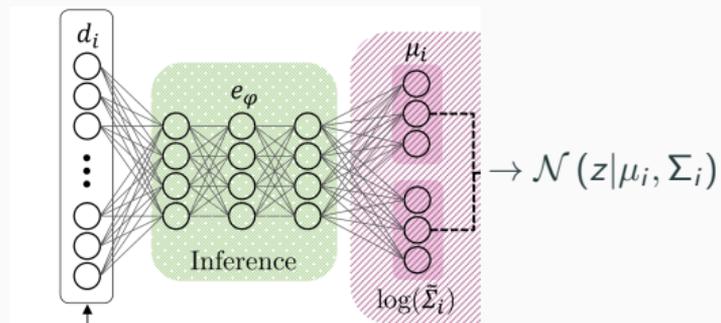
$$\mathcal{M}^{-1}(\mu) \leq \langle (z - \mu)(z - \mu)^{\dagger} \rangle_{\mathcal{P}(z|d)}$$



## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_\theta}{\partial \mathbf{z}} \right)^\dagger N^{-1} \left( \frac{\partial f_\theta}{\partial \mathbf{z}} \right) \right]_{\mathbf{z}=\mu} + \mathbf{1}$$

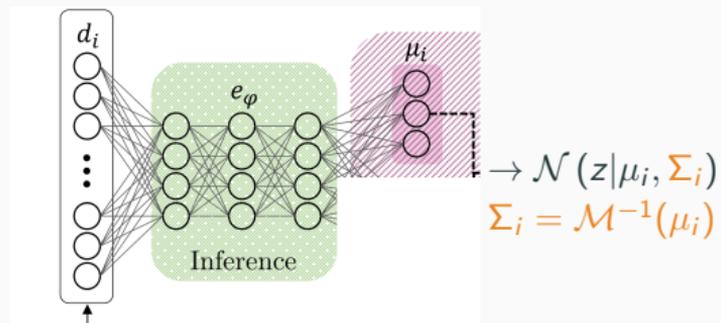
- † Kramer-Rao bound
$$\mathcal{M}^{-1}(\mu) \leq \langle (z - \mu)(z - \mu)^\dagger \rangle_{\mathcal{P}(z|d)}$$
- † Fisher-Metric is an approximation to post. covariance (including off-diagonal)



## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right)^{\dagger} N^{-1} \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right) \right]_{\mathbf{z}=\mu} + \mathbf{1}$$

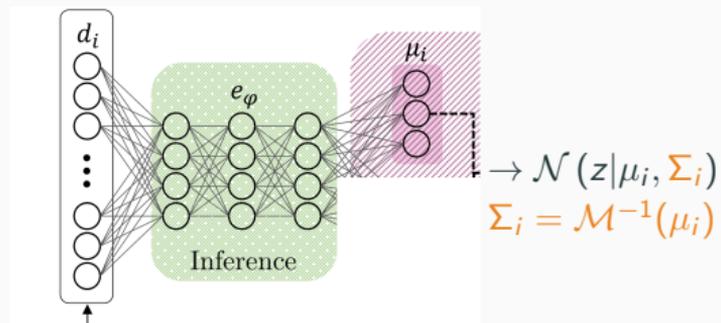
- † Kramer-Rao bound
$$\mathcal{M}^{-1}(\mu) \leq \langle (z - \mu)(z - \mu)^{\dagger} \rangle_{\mathcal{P}(z|d)}$$
- † Fisher-Metric is an approximation to post. covariance (including off-diagonal)



## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_\theta}{\partial z} \right)^\dagger N^{-1} \left( \frac{\partial f_\theta}{\partial z} \right) \right]_{z=\mu} + \mathbb{1}$$

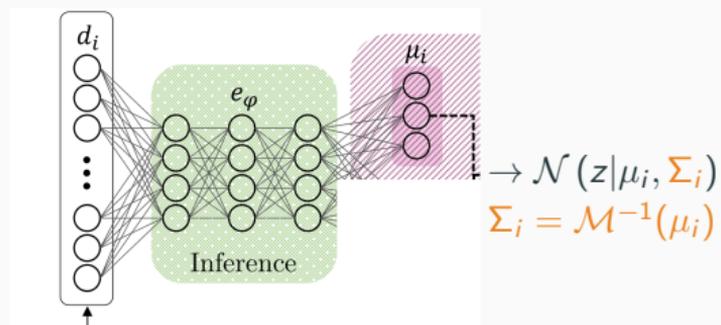
- † Kramer-Rao bound
$$\mathcal{M}^{-1}(\mu) \leq \langle (z - \mu)(z - \mu)^\dagger \rangle_{\mathcal{P}(z|d)}$$
- † Fisher-Metric is an approximation to post. covariance (including off-diagonal)
- † Coupled second moments of  $\mathcal{P}$  and  $\mathcal{Q}$



## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_{\theta}}{\partial z} \right)^{\dagger} N^{-1} \left( \frac{\partial f_{\theta}}{\partial z} \right) \right]_{z=\mu} + \mathbb{1}$$

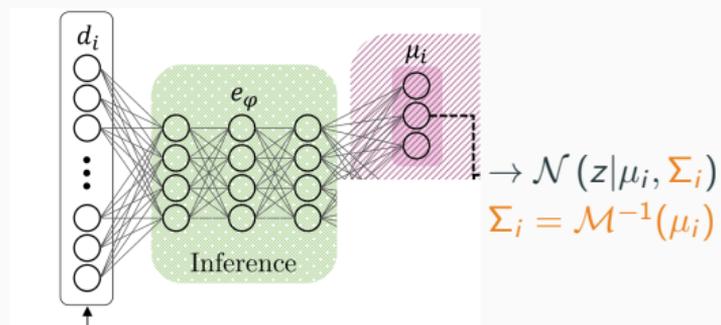
- † Kramer-Rao bound
$$\mathcal{M}^{-1}(\mu) \leq \langle (z - \mu)(z - \mu)^{\dagger} \rangle_{\mathcal{P}(z|d)}$$
- † Fisher-Metric is an approximation to post. covariance (including off-diagonal)
- † Coupled second moments of  $\mathcal{P}$  and  $\mathcal{Q}$
- † Increased inference capabilities of  $\mathcal{Q}$



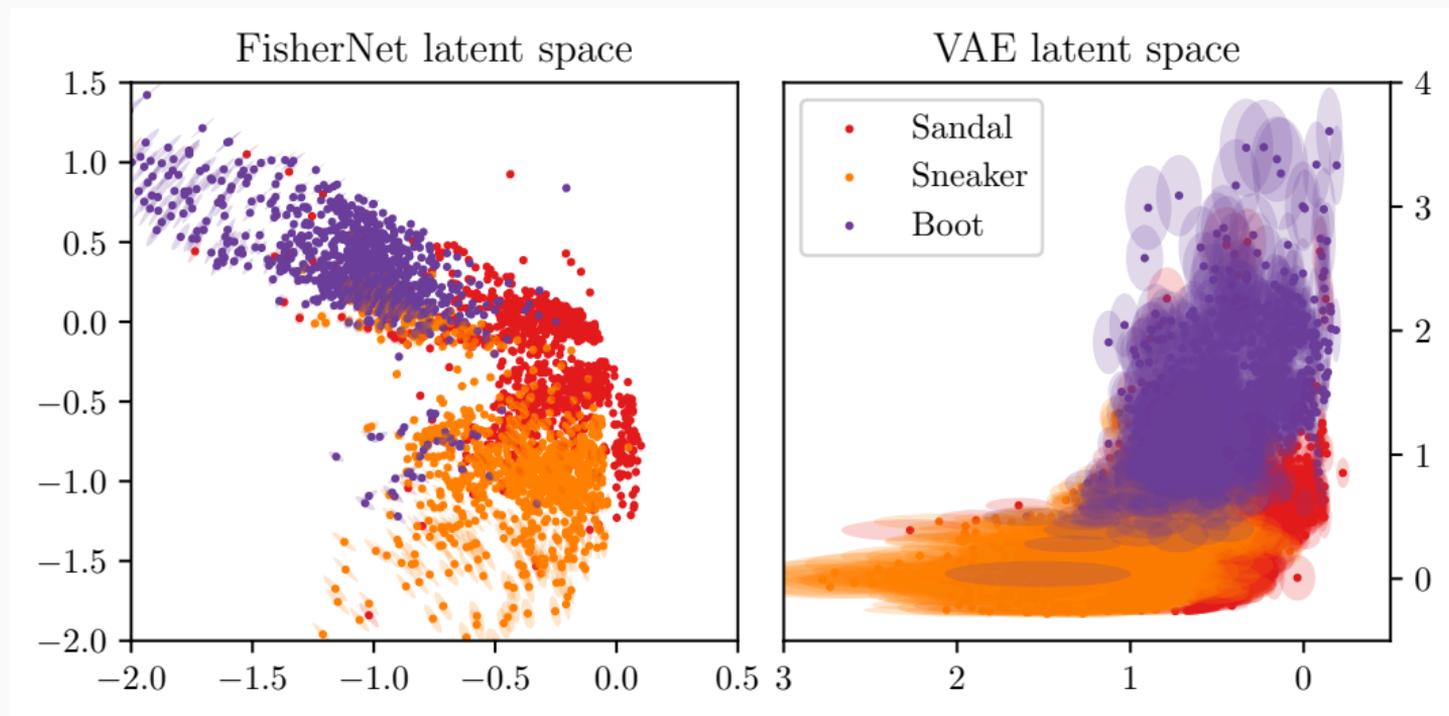
## Fisher Information metric

$$\mathcal{M}(\mu) \equiv \left[ \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right)^{\dagger} N^{-1} \left( \frac{\partial f_{\theta}}{\partial \mathbf{z}} \right) \right]_{\mathbf{z}=\mu} + \mathbb{1}$$

- † Kramer-Rao bound
$$\mathcal{M}^{-1}(\mu) \leq \langle (z - \mu)(z - \mu)^{\dagger} \rangle_{\mathcal{P}(z|d)}$$
- † Fisher-Metric is an approximation to post. covariance (including off-diagonal)
- † Coupled second moments of  $\mathcal{P}$  and  $\mathcal{Q}$
- † Increased inference capabilities of  $\mathcal{Q}$
- † Sampling is less trivial



## Comparison - Fashion MNIST



- ✦ Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...
- ✦ Change approximation distribution  $Q$ 
  - Non-Gaussian approximation
  - Use full covariance matrix
  - Estimate covariance from generative model using Fisher Information metric
  - ...

- ✦ Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...
- ✦ Change approximation distribution  $Q$ 
  - Non-Gaussian approximation
  - Use full covariance matrix
  - Estimate covariance from generative model using Fisher Information metric
  - ...
- ✦ Post-process encoded distribution  $Q$  in latent space

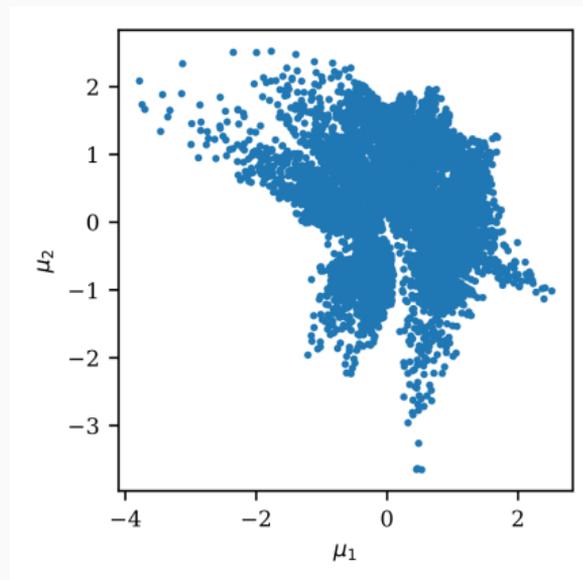
- ✦ Exchange / Modify likelihood distribution
  - Noise covariance  $N$  becomes a variable
  - ...
- ✦ Change approximation distribution  $Q$ 
  - Non-Gaussian approximation
  - Use full covariance matrix
  - Estimate covariance from generative model using Fisher Information metric
  - ...
- ✦ Post-process encoded distribution  $Q$  in latent space
  - Transform latent space distribution to a standard normal distribution for sample generation
  - ...

## Transform latent space

† Marginal approximate prior distribution  $Q(z) = \int Q(z|d) P(d) \, dd$

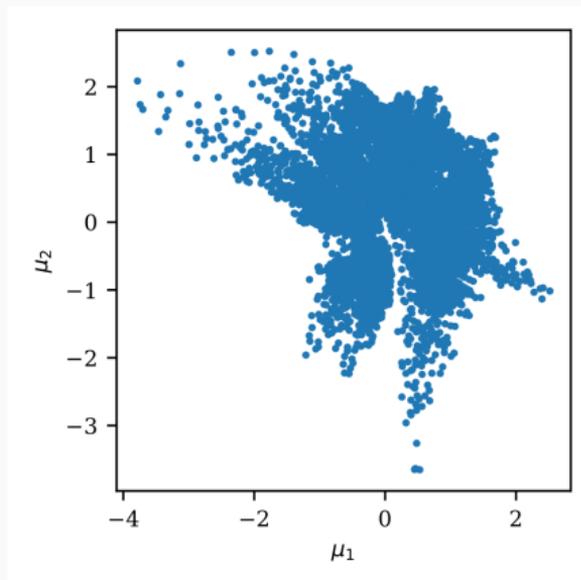
## Transform latent space

† Marginal approximate prior distribution  $Q(z) = \int Q(z|d) P(d) dd$



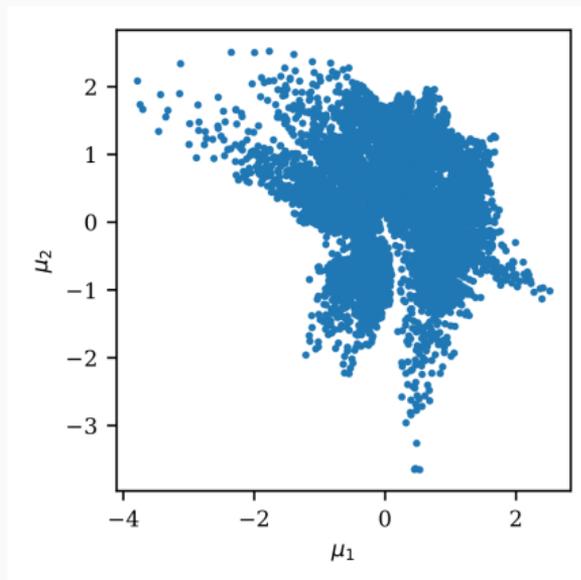
## Transform latent space

- + Marginal approximate prior distribution  $Q(z) = \int Q(z|d) P(d) dd$
- + Ideally  $Q(z) \simeq \mathcal{P}(z) = \mathcal{N}(z|0, \mathbf{1})$

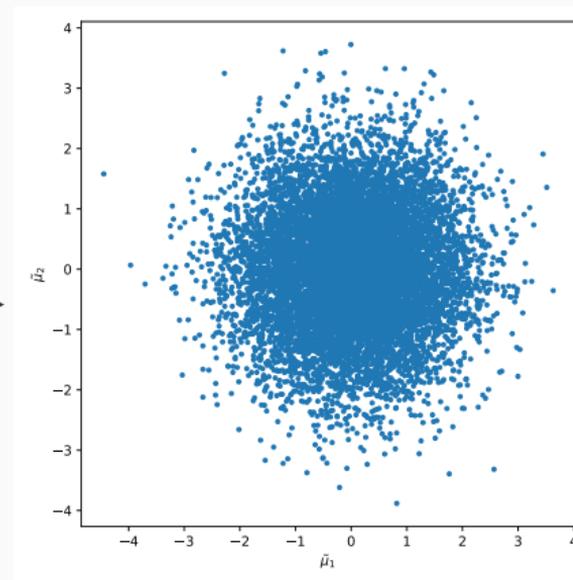


# Transform latent space

- + Marginal approximate prior distribution  $Q(z) = \int Q(z|d) P(d) dd$
- + Ideally  $Q(z) \simeq \mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$

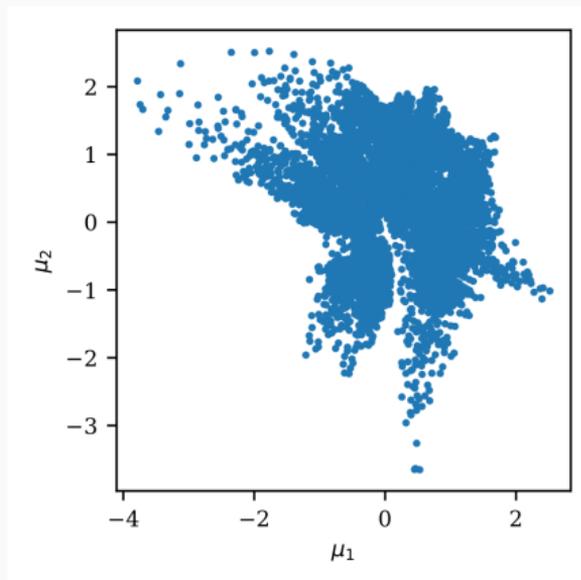


Transform using e.g.  
Normalizing flows or  
density estimation



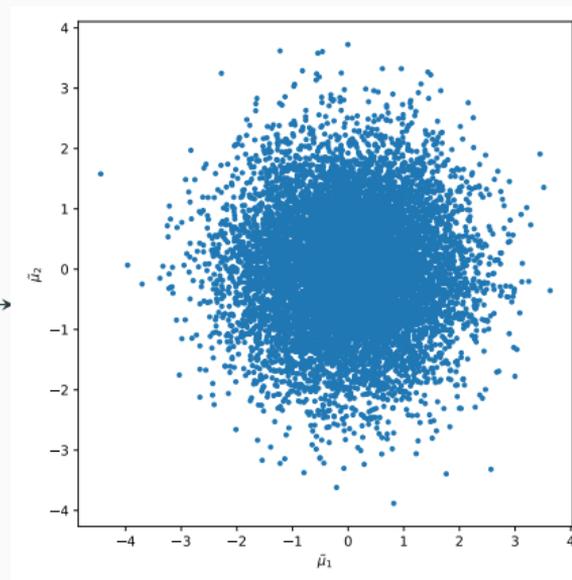
# Transform latent space

- + Marginal approximate prior distribution  $Q(z) = \int Q(z|d) P(d) dd$
- + Ideally  $Q(z) \simeq \mathcal{P}(z) = \mathcal{N}(z|0, \mathbb{1})$



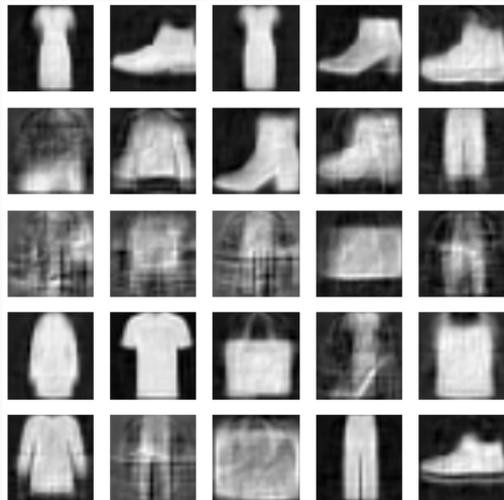
Transform using e.g.  
Normalizing flows or  
density estimation

←  
Use inverse trans-  
formation for sample  
generation



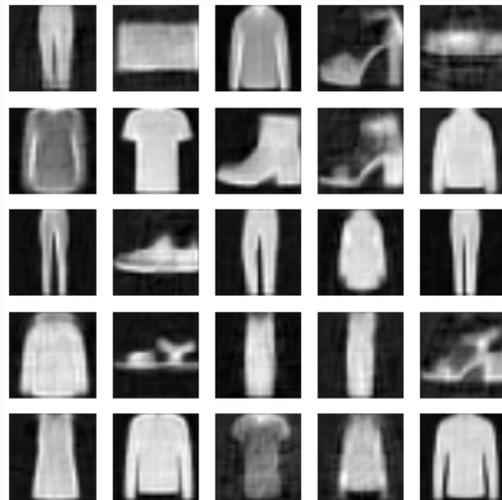
# Transform latent space

Drawn from latent space



FID: 110

Drawn from transformed space



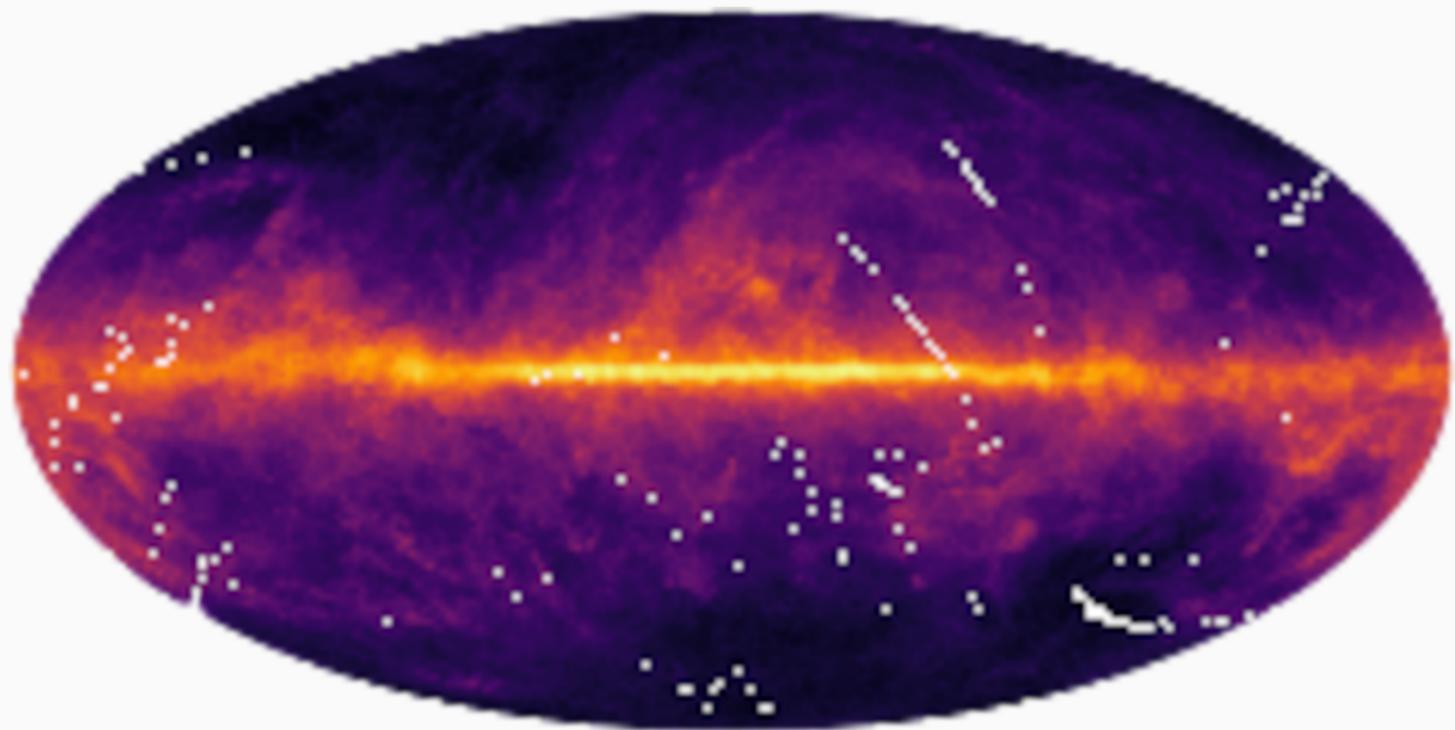
(Training FID: 80)

FID: 85

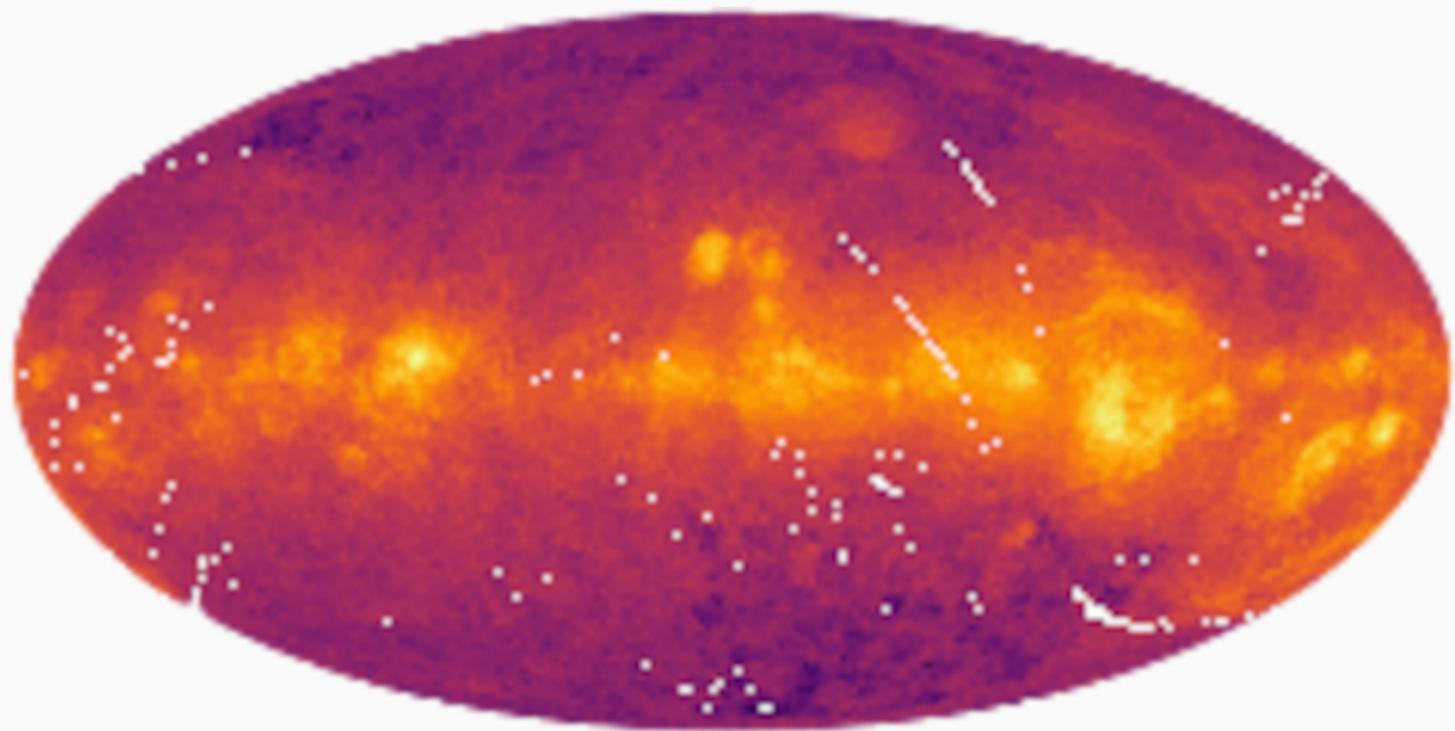
# Decomposition of the Galactic multi-frequency sky

---

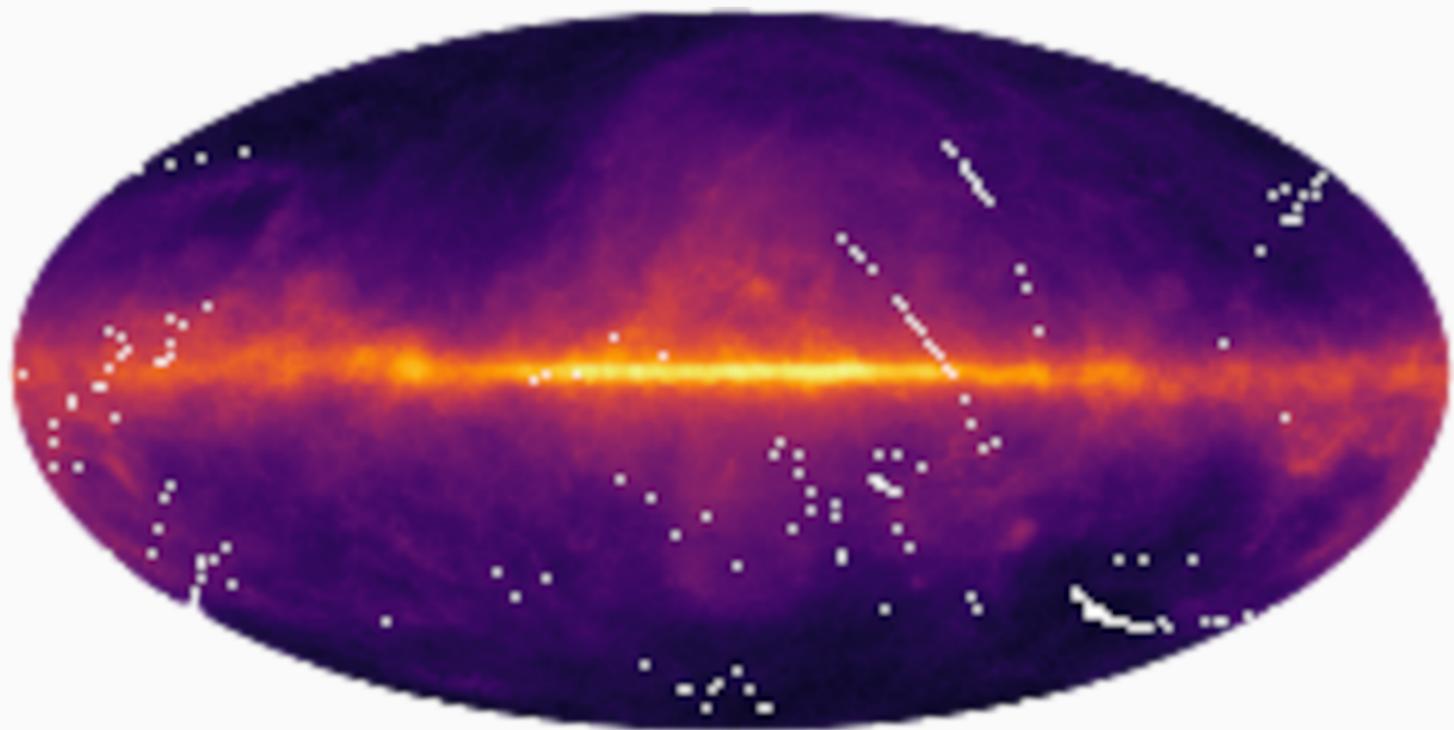
## Planck Map (545 GHz)



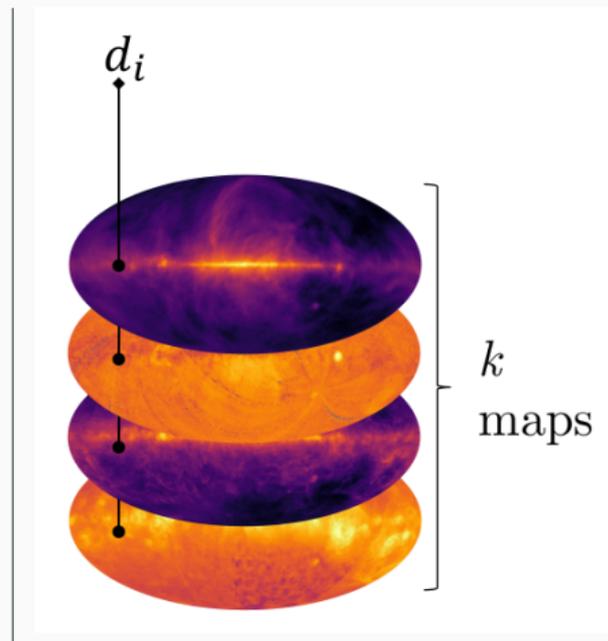
H-alpha (656.3 nm)



# Fermi $\gamma$ -ray map (1.7 GeV)

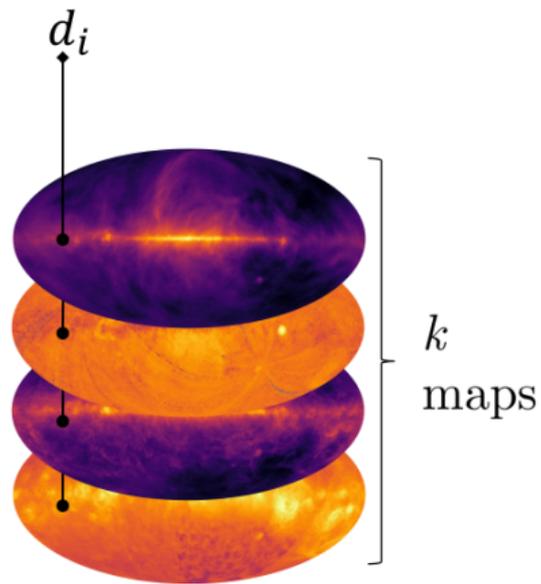


# Decomposition of the Galactic multi-frequency sky - Data



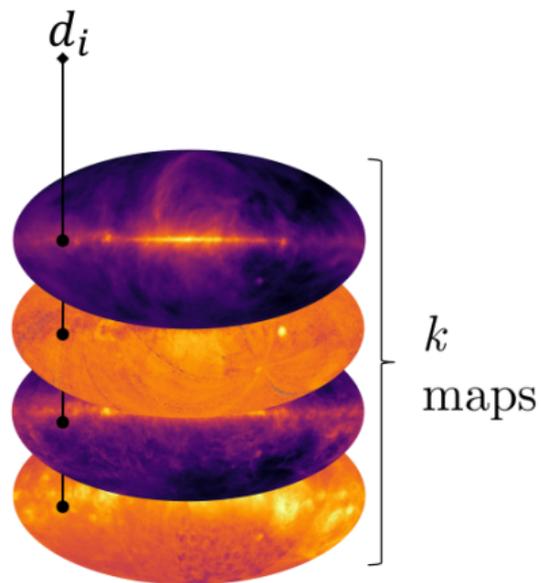
# Decomposition of the Galactic multi-frequency sky - Data

- + Pixel based approach ( $d_i$ : frequency brightness vector at  $i$ th-location on the sky)



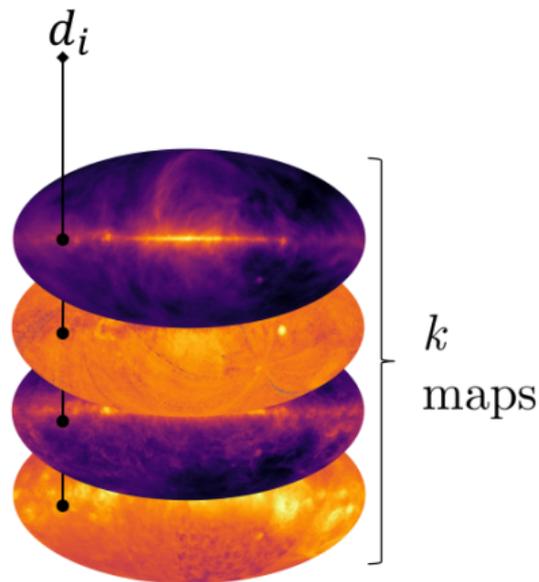
## Decomposition of the Galactic multi-frequency sky - Data

- + Pixel based approach ( $d_i$ : frequency brightness vector at  $i$ th-location on the sky)
- +  $k = 35$  all-sky maps



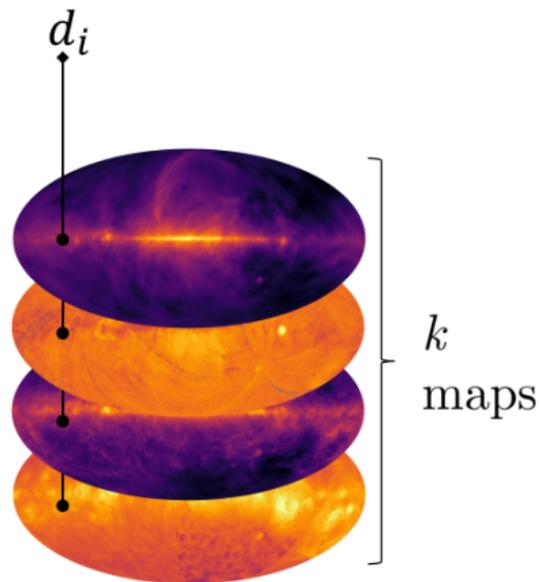
# Decomposition of the Galactic multi-frequency sky - Data

- + Pixel based approach ( $d_i$ : frequency brightness vector at  $i$ th-location on the sky)
- +  $k = 35$  all-sky maps
- + Frequencies from Radio (MHz) to gamma ray (GeV)

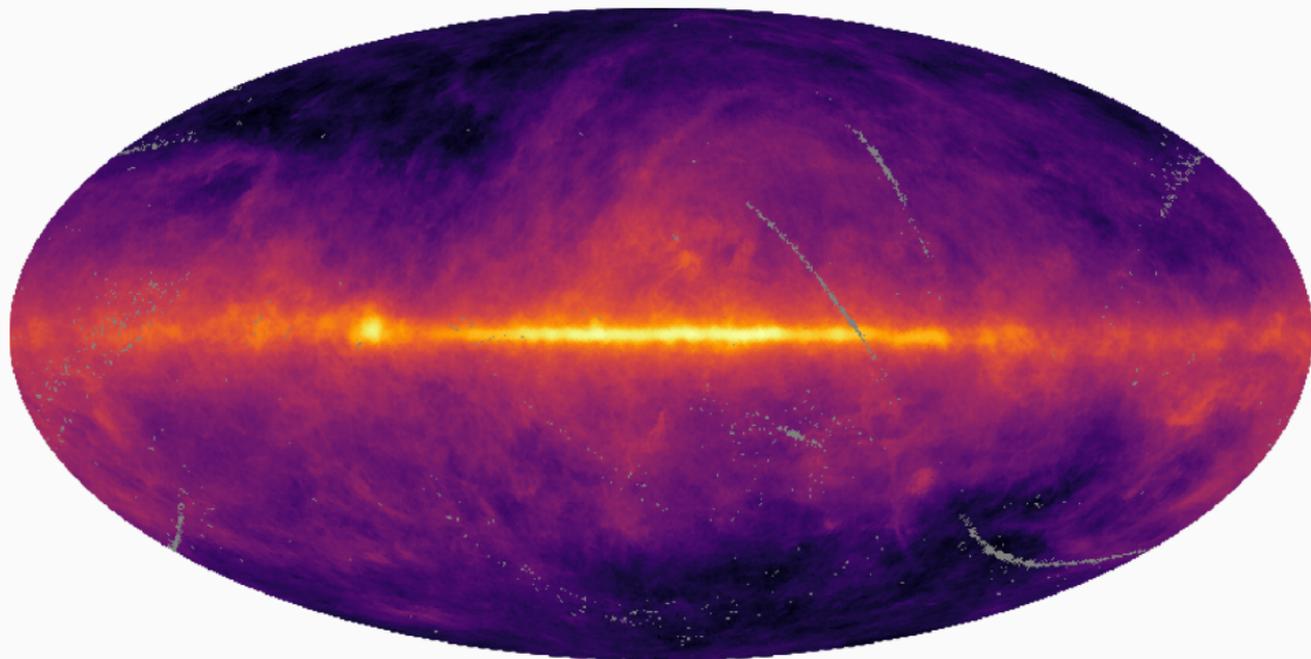


# Decomposition of the Galactic multi-frequency sky - Data

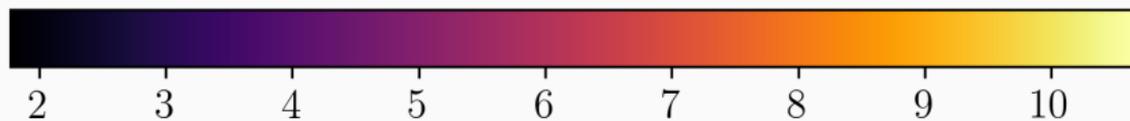
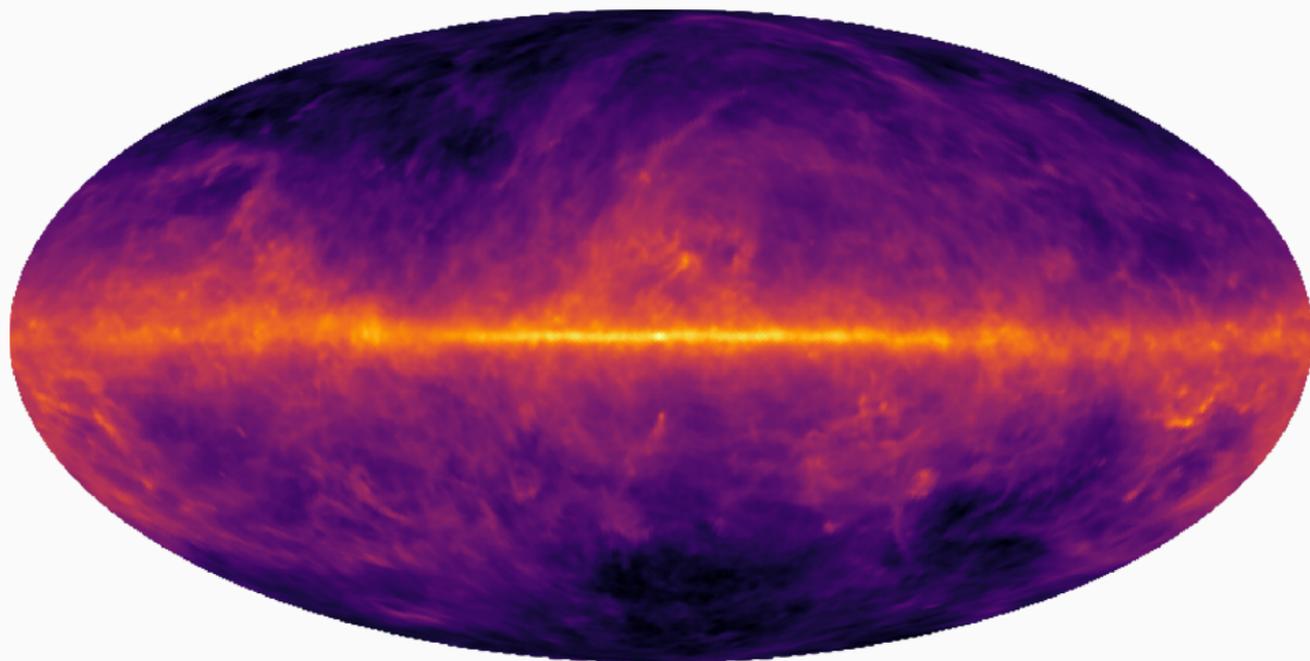
- + Pixel based approach ( $d_i$ : frequency brightness vector at  $i$ th-location on the sky)
- +  $k = 35$  all-sky maps
- + Frequencies from Radio (MHz) to gamma ray (GeV)
- + 10 dimensional latent space



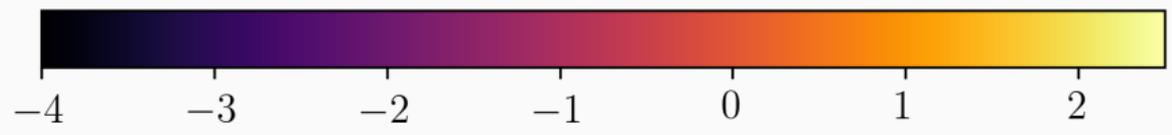
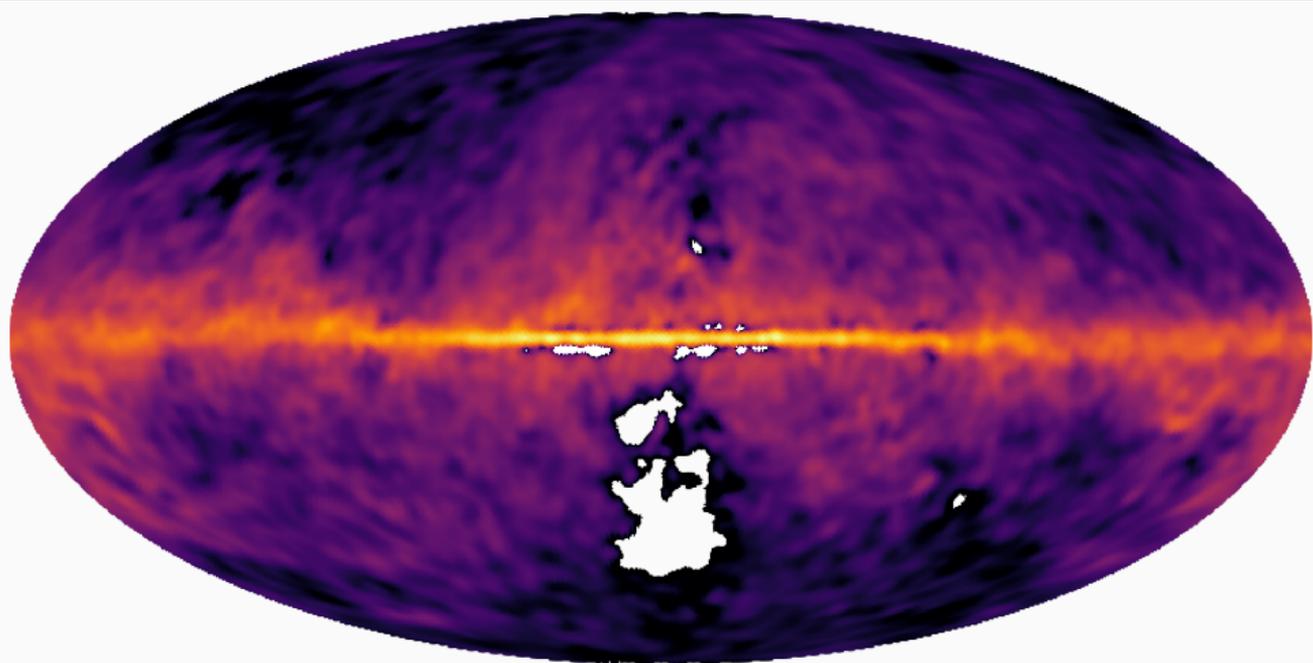
## Decomposition of the Galactic multi-frequency sky - Feature A



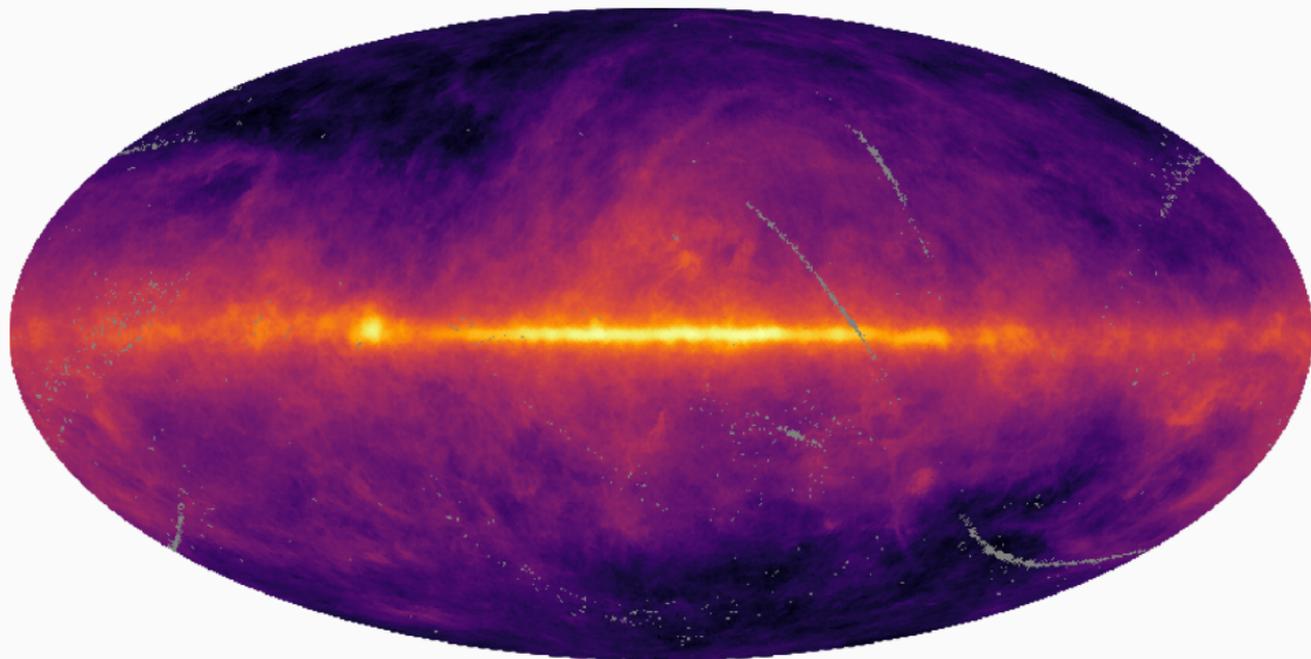
## Decomposition of the Galactic multi-frequency sky - Dust



# Decomposition of the Galactic multi-frequency sky - Hadronic Component

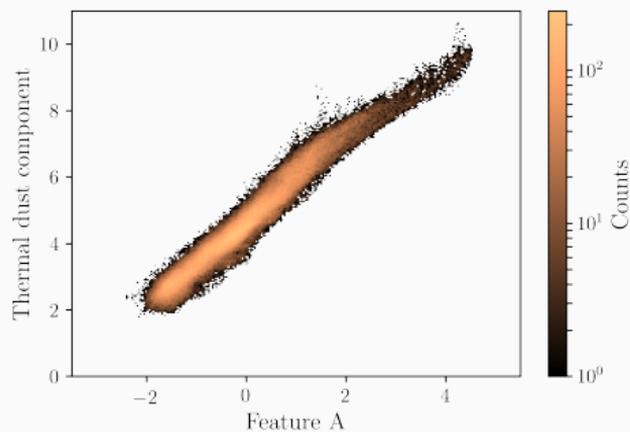


## Decomposition of the Galactic multi-frequency sky - Feature A



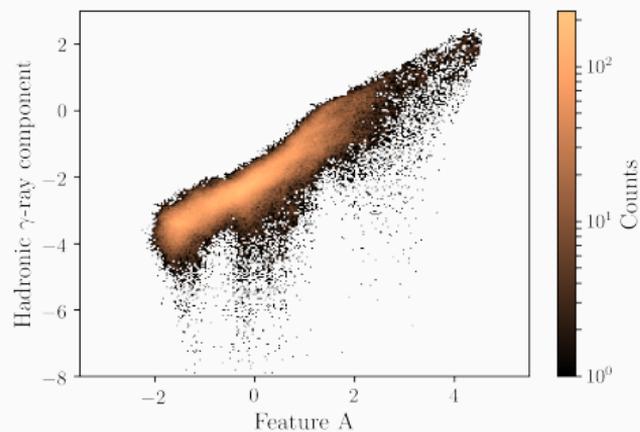
# Decomposition of the Galactic multi-frequency sky - Correlations

Correlation Dust



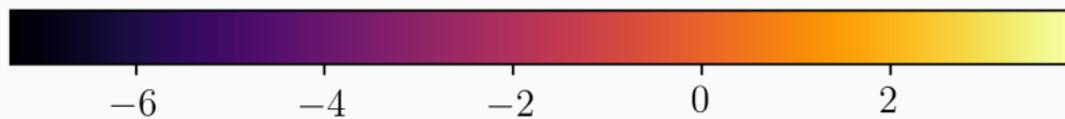
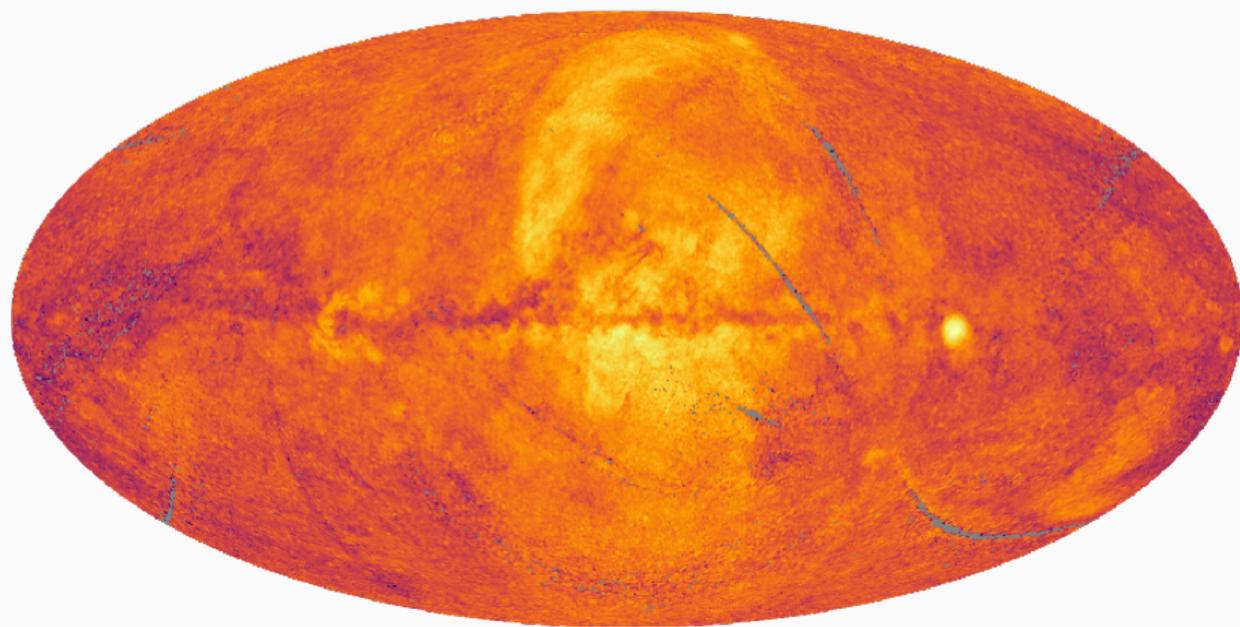
Mutual Information:  $I(X; Y) = 1.72$

Correlation Hadronic Component



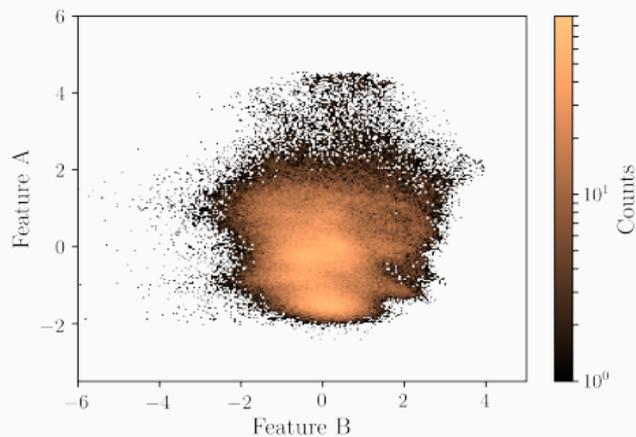
Mutual Information:  $I(X; Y) = 1.07$

# Decomposition of the Galactic multi-frequency sky - Feature B



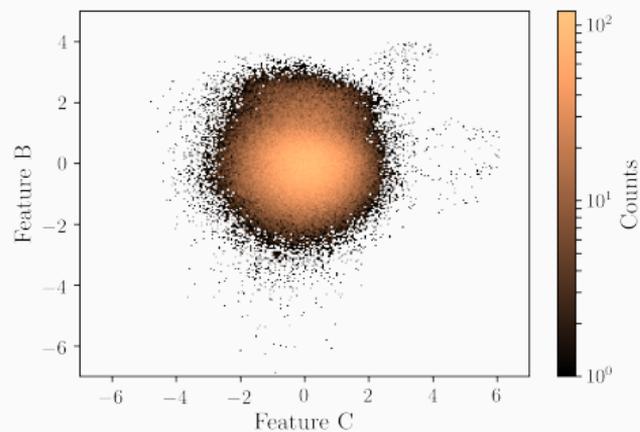
# Decomposition of the Galactic multi-frequency sky - Correlations

Correlation Feature A



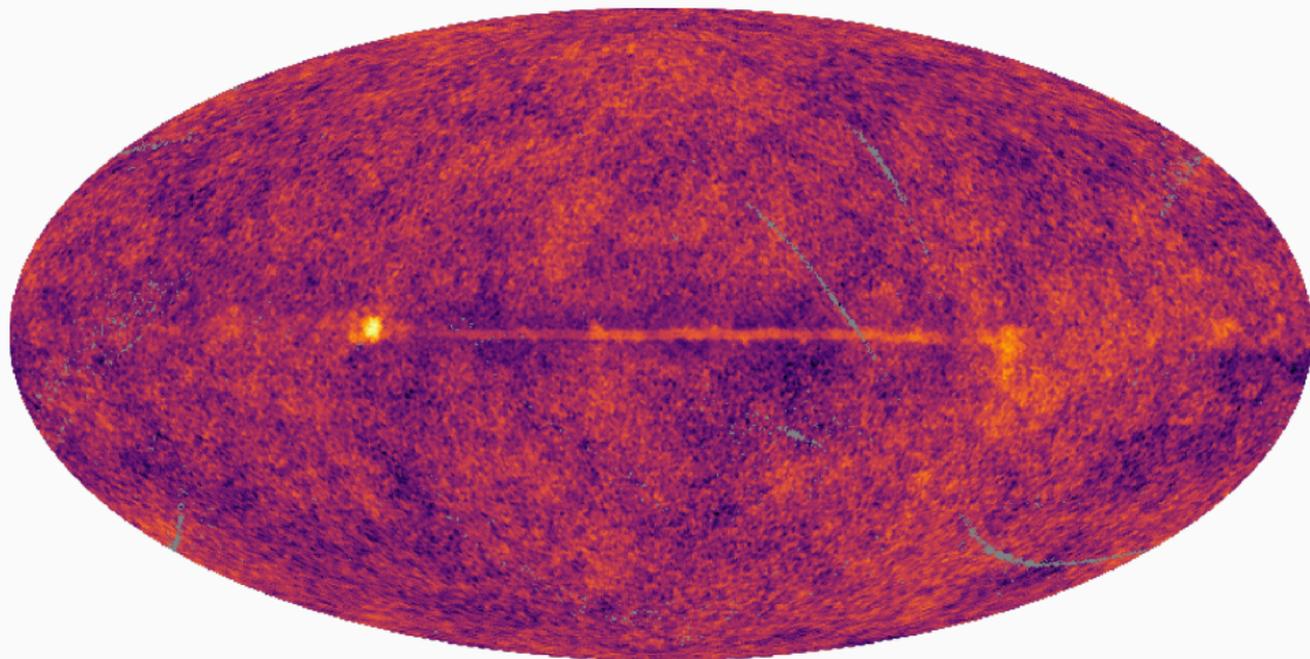
Mutual Information:  $I(X; Y) = 0.33$

Correlation Feature C

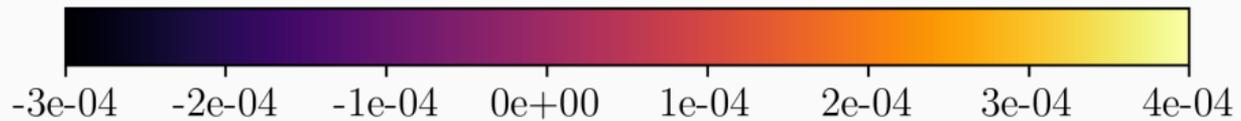
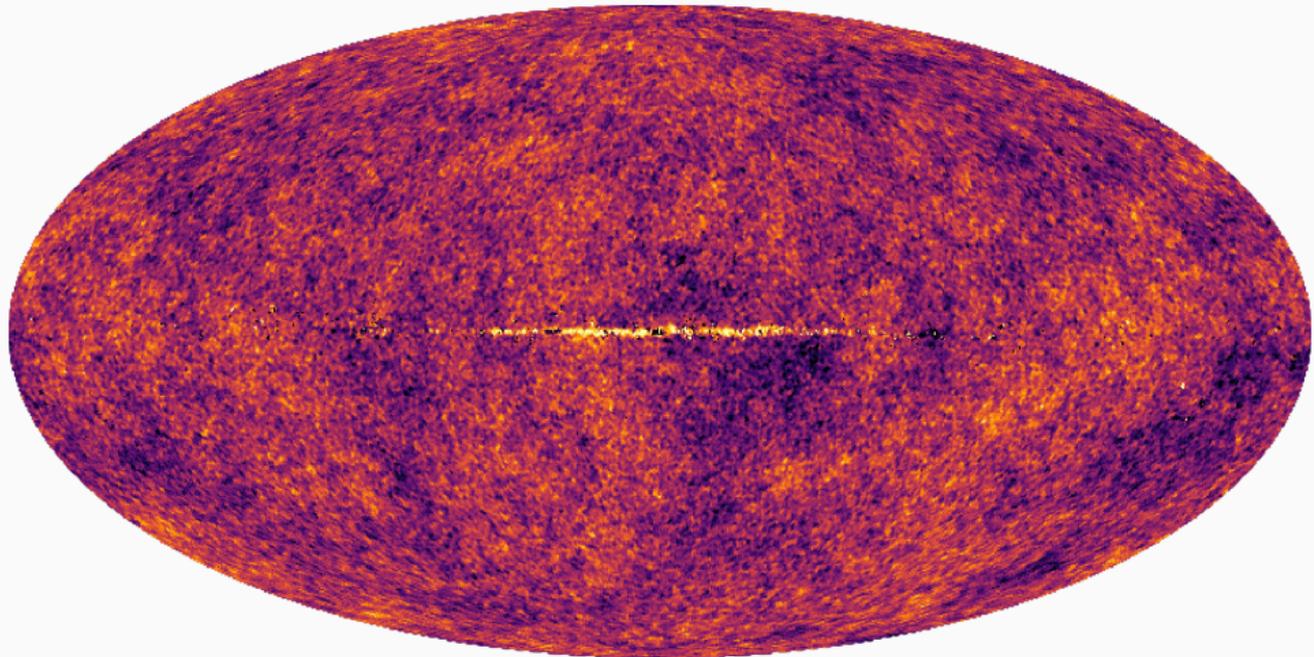


Mutual Information:  $I(X; Y) = 0.20$

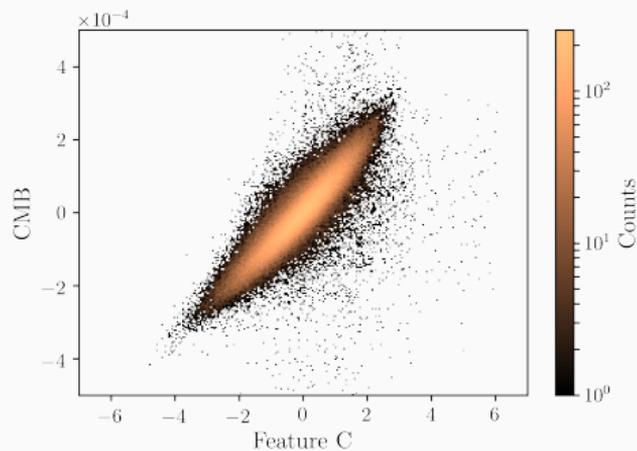
# Decomposition of the Galactic multi-frequency sky - Feature C



# Decomposition of the Galactic multi-frequency sky - CMB



Correlation CMB



Mutual Information:  $I(X; Y) = 0.88$

- + VAEs use probabilistic generative models

- + VAEs use probabilistic generative models
- + Associated inverse problem solved using variational inference

- + VAEs use probabilistic generative models
- + Associated inverse problem solved using variational inference
- + Fisher information metric as covariance can improve inference capacity
- + Posterior analysis of latent space using normalizing transformations can improve sampling quality

- + VAEs use probabilistic generative models
- + Associated inverse problem solved using variational inference
- + Fisher information metric as covariance can improve inference capacity
- + Posterior analysis of latent space using normalizing transformations can improve sampling quality
- + Latent features of multi-frequency sky partially coincide with known physical components superimposed on the sky

- † **Bayesian decomposition of the Galactic multi-frequency sky using probabilistic autoencoders**; Milosevic, S., Frank, P., Leike, R., Müller, A., Enßlin, T.A.; *Astronomy & Astrophysics* 2021, 650, A100. doi.org/10.1051/0004-6361/202039435
- † **Probabilistic Autoencoder using Fisher Information**; Zacherl, J., Frank, P., and Enßlin, T.A.; *Entropy* 2021, 23, 1640. doi.org/10.3390/e23121640