

Massive Black Hole Binary parameter estimation using Masked Autoregressive Flows

Ivan Martin Vilchez

LISA Data Analysis: from Classical Methods to Machine Learning

L2IT Toulouse, November 21–25 2022

Institute of
Space Sciences



CSIC

CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS



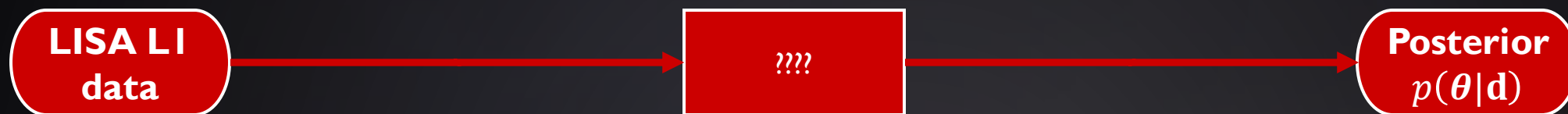
IEEC[®]

Institut d'Estudis
Espacials de Catalunya

UAB

Universitat Autònoma
de Barcelona

Towards a LISA DA pipeline



Bayesian inference

Likelihood $\log p(\mathbf{d} | \boldsymbol{\theta}) = \sum_{n=\{A,E,T\}} 2 (\mathbf{h}(\boldsymbol{\theta}) | \mathbf{d})_n - (\mathbf{h}(\boldsymbol{\theta}) | \mathbf{h}(\boldsymbol{\theta}))_n$

forward simulations! integrals!

$$p(\boldsymbol{\theta} | \mathbf{d}) = \frac{p(\mathbf{d} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{d})}$$

Posterior ← Prior

Evidence even harder integrals!

$$p(\mathbf{d}) = \int p(\mathbf{d} | \boldsymbol{\theta}) d\boldsymbol{\theta}$$

Density estimation

- In practice, “classical” methods usually sample $p(\boldsymbol{\theta} \mid \mathbf{x}_{\text{obs}})$ with e.g. MCMC.
 - Use $p(\boldsymbol{\theta} \mid \mathbf{x}_{\text{obs}}) \propto p(\mathbf{x}_{\text{obs}} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta})$ to get posterior samples.
 - Sampling the likelihood can be tricky and take a long time...
- We'll try avoid computing $p(\mathbf{x} \mid \boldsymbol{\theta})$ at all using Machine Learning.
Options:
 1. Estimate $p(\mathbf{x}, \boldsymbol{\theta})$, then substitute \mathbf{x}_{obs} to find $p(\boldsymbol{\theta} \mid \mathbf{x}_{\text{obs}})$
 2. Estimate $p(\boldsymbol{\theta} \mid \mathbf{x})$ directly, then substitute \mathbf{x}_{obs}
 3. Estimate $p(\mathbf{x} \mid \boldsymbol{\theta})$ and use Bayes to find $p(\boldsymbol{\theta} \mid \mathbf{x}_{\text{obs}})$

Normalising Flows

Consider a complicated “target distribution” $q(x)$.
Can we find an invertible transformation $x = f(z)$ such that, starting from a simple “base distribution” $p(z)$, we recover $q(x)$ via change of variable?

$$\int q(x) dx = \int p(z) dz = 1$$

$$q(x) = p(z) \left| \frac{dz}{dx} \right| = p(f^{-1}(x)) \left| \frac{df^{-1}(x)}{dx} \right|$$

For multivariate distributions,

$$q(\mathbf{x}) = p(\mathbf{z}) |\det \mathbf{J}[f^{-1}(x)]|$$

↖
In practice, this determinant is generally **expensive** to compute!

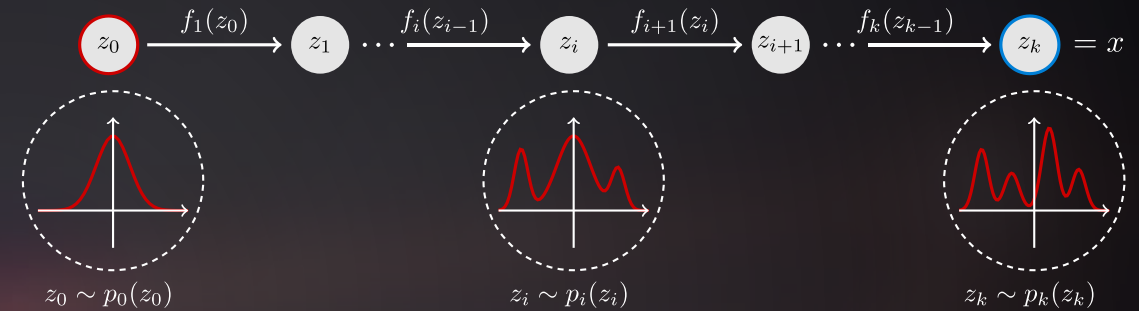


Image: [Janosh Riebesell @ github](#)

To create complex transformations, use compositions of simpler ones → a *flow* of transformations:

$$f(\mathbf{z}) = [f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1](\mathbf{z})$$

$$q(\mathbf{x}) = p(\mathbf{z}) \prod_{i=1}^k |\det \mathbf{J}[f_i(\mathbf{z}_{i-1})]|^{-1}$$

Masked Autoregressive Flows (MAF)

In summary: we will be chaining some transformation $f_i(\mathbf{z})$, which needs to be:

- Invertible
- Differentiable
- Easy-to-compute $\det \mathbf{J}$ (e.g., triangular \mathbf{J})
- Expressive!

Several exist in the literature.

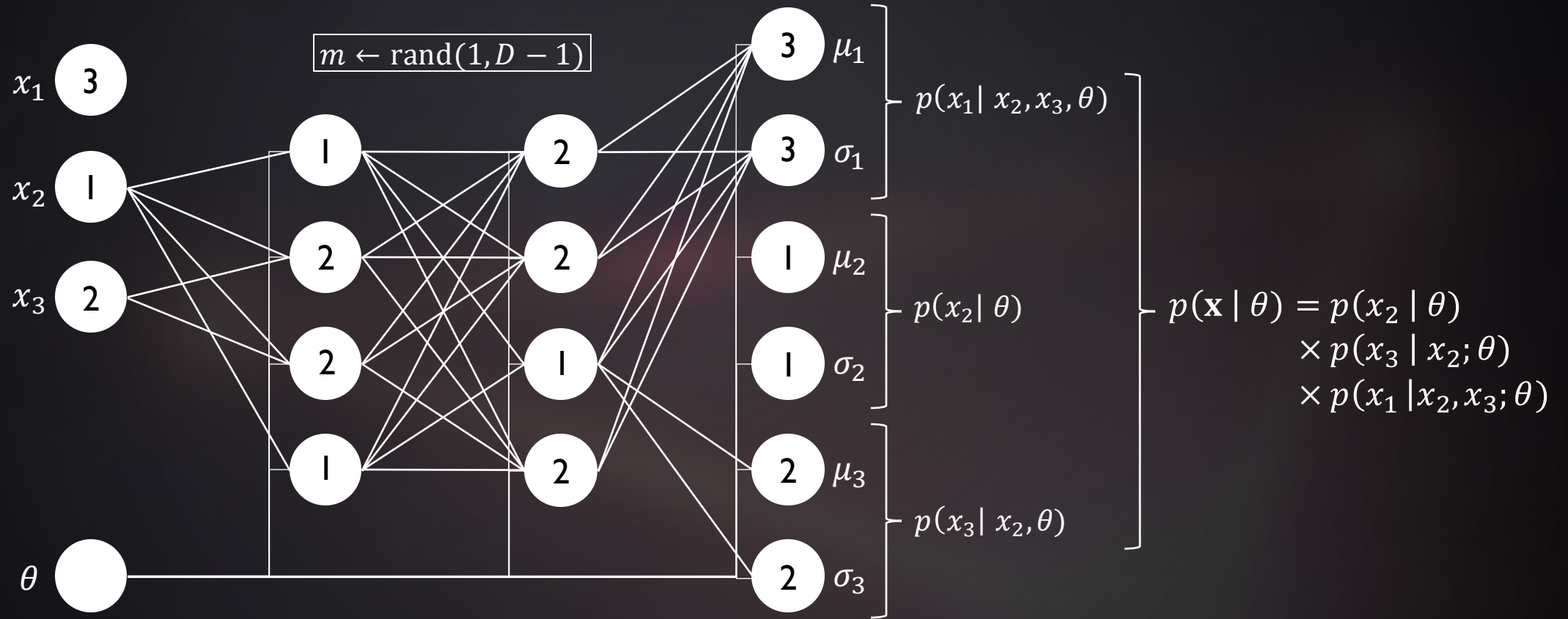
We choose: $x_j(\mathbf{z}) = \mu_j + z_j \sigma_j$ $\begin{cases} \mu_j = f_{\mu_j}(\mathbf{x}_{1,\dots,j-1}), \\ \sigma_j = f_{\sigma_j}(\mathbf{x}_{1,\dots,j-1}) \end{cases}$

This is equivalent to modeling

$$q(\mathbf{x}) = \prod_{j=1}^D p(x_j | \mathbf{x}_{1,\dots,j-1}) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \dots \rightarrow \text{This is our loss function!}$$

\swarrow
 $\mathcal{N}(x_j | \mu_j(x_{1,\dots,j-1}), \sigma_j^2(x_{1,\dots,j-1}))$

Masked Autoencoder for Distribution Estimation (MADE)



MADE in practice: Masks

It is computationally easier to introduce masks \mathbf{M} so that:

$$\mathbf{h}^l(\mathbf{h}^{l-1}) = \mathbf{g}(\mathbf{W}^l \odot \mathbf{M}^{\mathbf{W}^l} \mathbf{h}^{l-1} + \mathbf{b}^l + \mathbf{W}_\theta^l \boldsymbol{\theta})$$

For hidden layers:

$$\mathbf{M}_{k',k}^{\mathbf{W}^l} = \begin{cases} 1 & \text{if } m^l(k') \geq m^{l-1}(k) \\ 0 & \text{otherwise.} \end{cases}$$

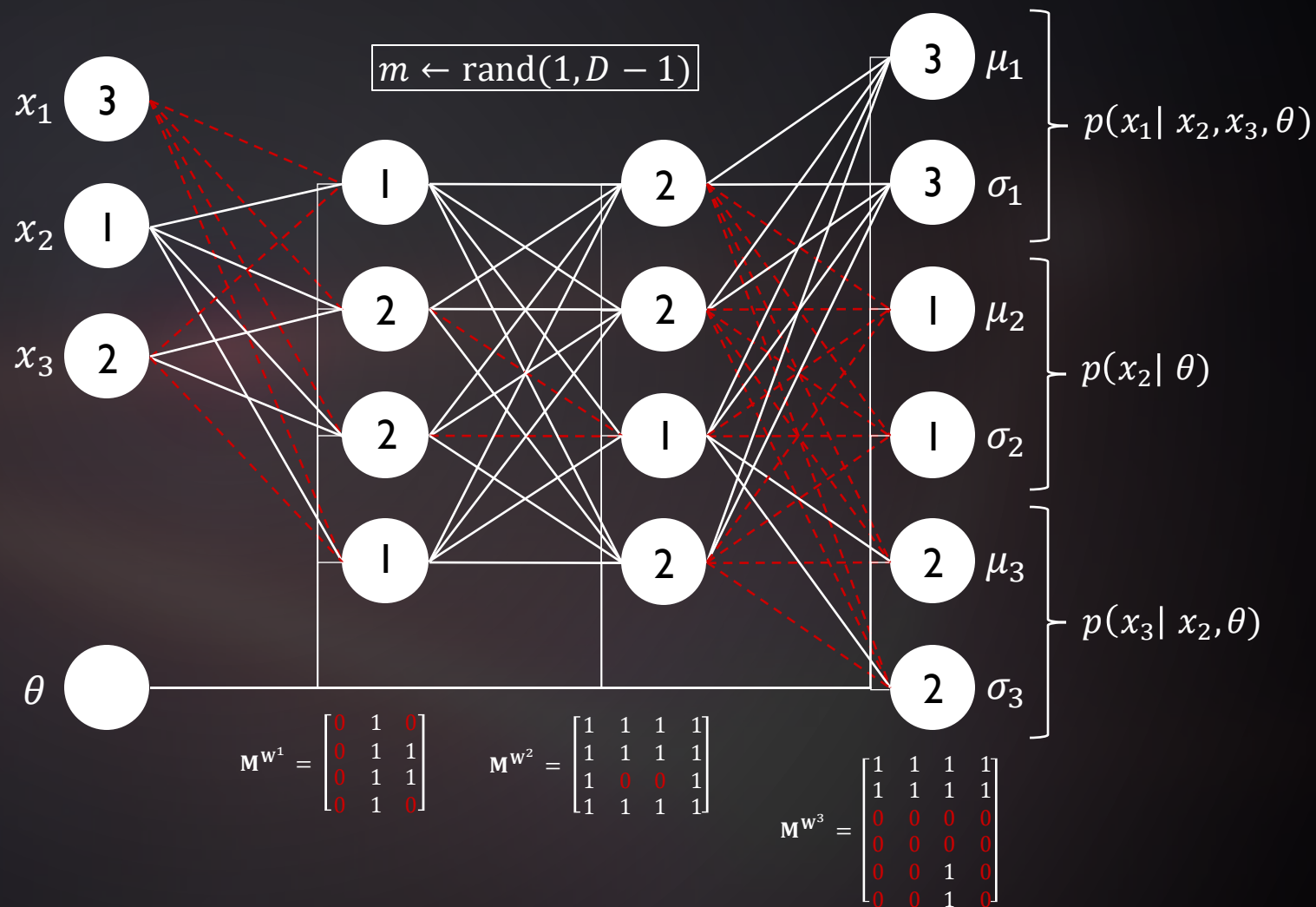
For the output layer:

$$\mathbf{M}_{k',k}^{\mathbf{W}^L} = \begin{cases} 1 & \text{if } m^L(k') > m^{L-1}(k) \\ 0 & \text{otherwise.} \end{cases}$$

Stack several of these,
with randomized ordering!



MAF



Towards a LISA DA pipeline

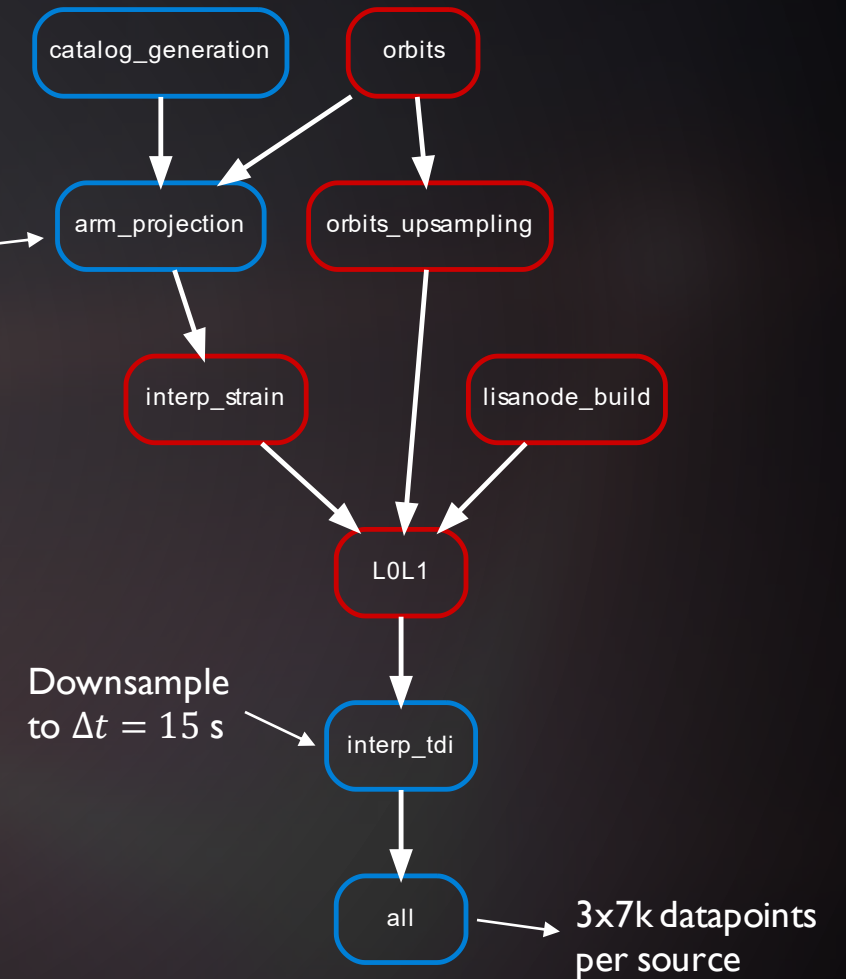


Data generation



20k(+2k) sources, sampled from $p(\theta)$

1 day around merger



Principal Component Analysis (PCA)

▪ Define

$$X = \left[\begin{array}{c|c|c|c} | & | & \dots & | \\ \mathbf{d}_1 & \mathbf{d}_2 & \dots & \mathbf{d}_n \\ | & | & \dots & | \end{array} \right] \left. \vphantom{\begin{array}{c|c|c|c} | & | & \dots & | \\ \mathbf{d}_1 & \mathbf{d}_2 & \dots & \mathbf{d}_n \\ | & | & \dots & | \end{array}} \right\} \begin{array}{l} (m \text{ rows}) \\ \\ \\ (n \text{ columns}) \end{array} \quad (m \geq n)$$

▪ One can decompose $X = U \Sigma V^T \longrightarrow \mathbf{d}_j = U \Sigma V_j^T$

$\begin{array}{ccc} (m, n) & \begin{array}{ccc} (m, m) & (m, n) & (n, n) \\ \text{unitary} & \text{diag.} & \text{unitary} \end{array} & \longrightarrow & \text{Singular Value Decomposition (SVD)} \end{array}$

▪ Approximately, one can decompose

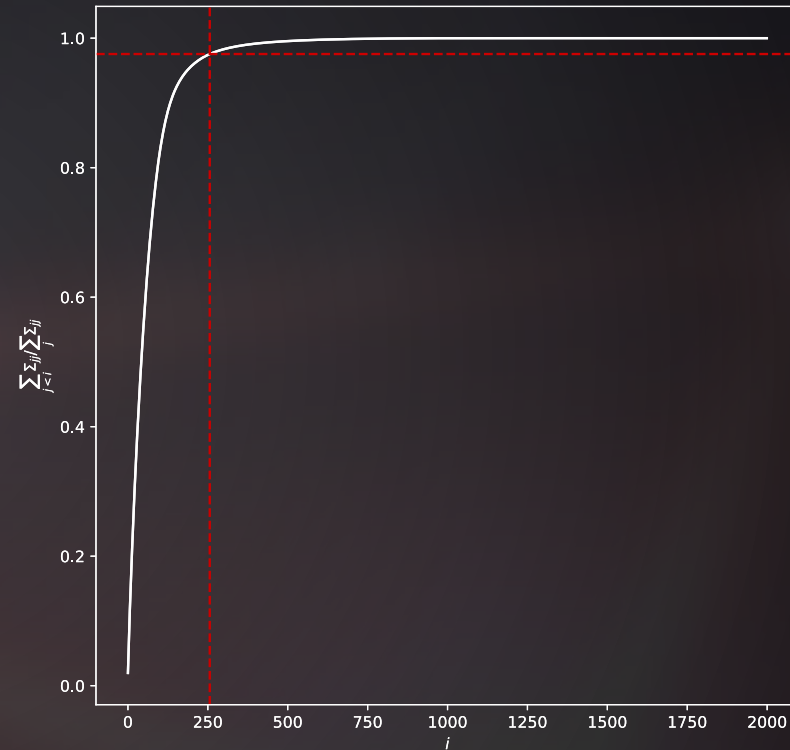
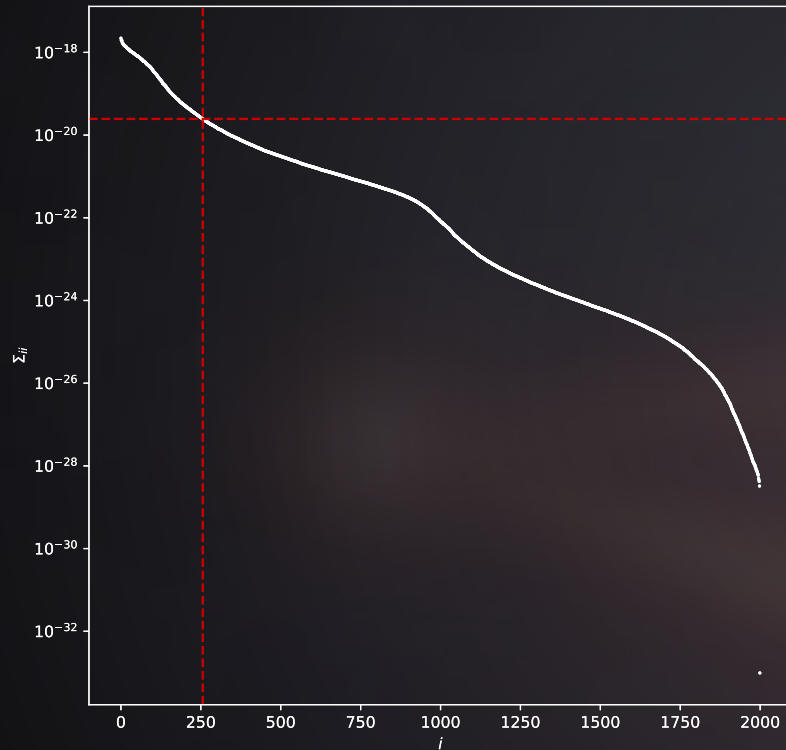
$$X \approx \tilde{X}_r = \tilde{U} \tilde{\Sigma} \tilde{V}^T \quad (r < n)$$

$\begin{array}{ccc} (m, n) & \begin{array}{ccc} (m, r) & (r, r) & (r, n) \end{array} \end{array}$

▪ Then, for new $\mathbf{d} \rightarrow \mathbf{x} = \tilde{U}^T \mathbf{d}$ is a data summary!

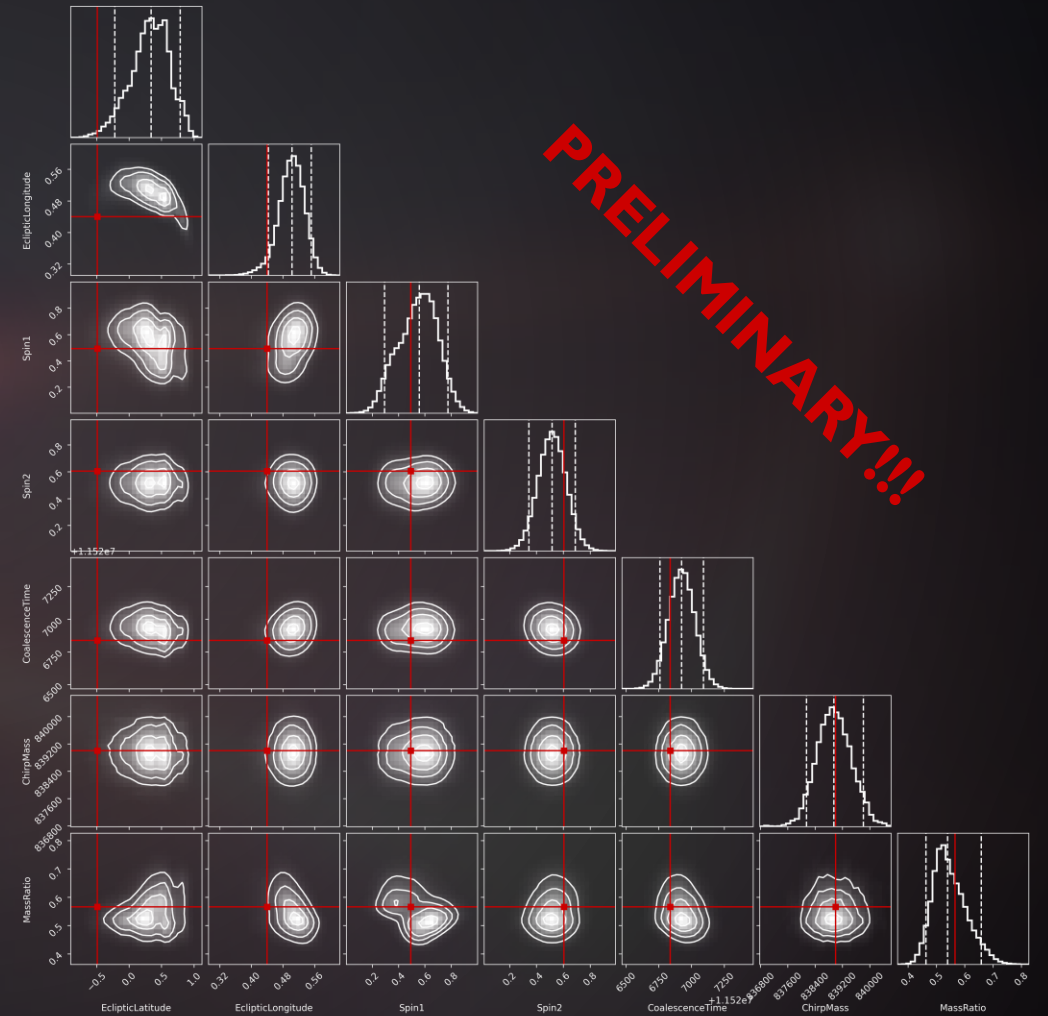
$\begin{array}{ccccc} (m) & \longrightarrow & (r) & \begin{array}{cc} (r, m) & (m) \end{array} \end{array}$

Principal Component Analysis (PCA)



Some results

Validation



Some results

Spritz



Some results

"Spritz-like"



Conclusions

- We were able to successfully model the log-likelihood as a MAF
- Applications:
 - Very fast evaluation of the likelihood
 - GPU-accelerated MCMC, with gradients!
 - Fast waveform emulators, possibly?
- Caveats:
 - Tuned to narrow prior (but should be generalizable with enough compute power)
 - Issues with actual LDC data
 - Plenty of fine-tuning avenues