

**MANGO**

# **Machine-Learning for Next-Generation GW Observatories**

Antsa Rasamoela

Research Engineer

LISA Consortium Core Member





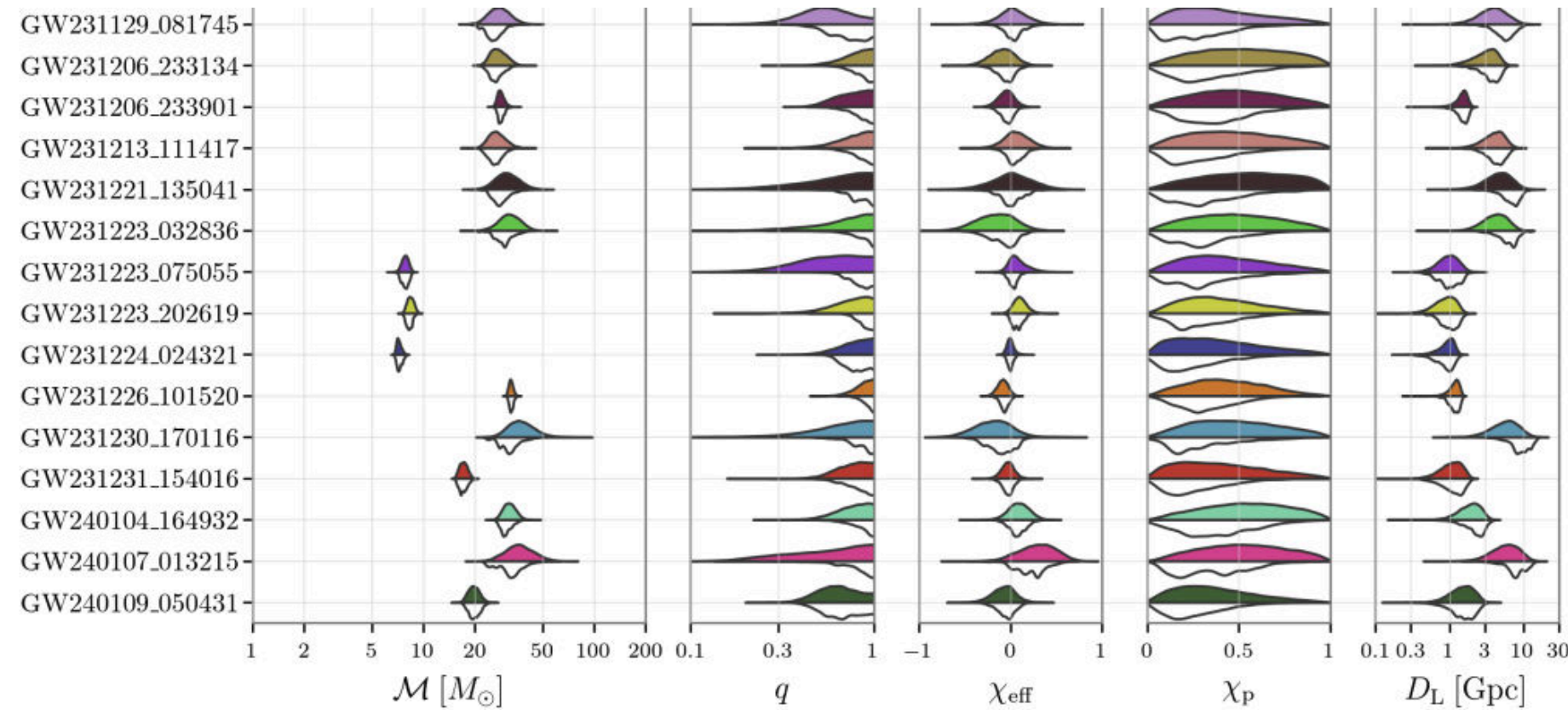
# GWTC-4.0: Updating the Gravitational-Wave Transient Catalog with Observations from the First Part of the Fourth LIGO-Virgo-KAGRA Observing Run

THE LIGO SCIENTIFIC COLLABORATION, THE VIRGO COLLABORATION, AND THE KAGRA COLLABORATION  
(SEE THE END MATTER FOR THE FULL LIST OF AUTHORS)

(Compiled: September 7, 2025)

## ABSTRACT

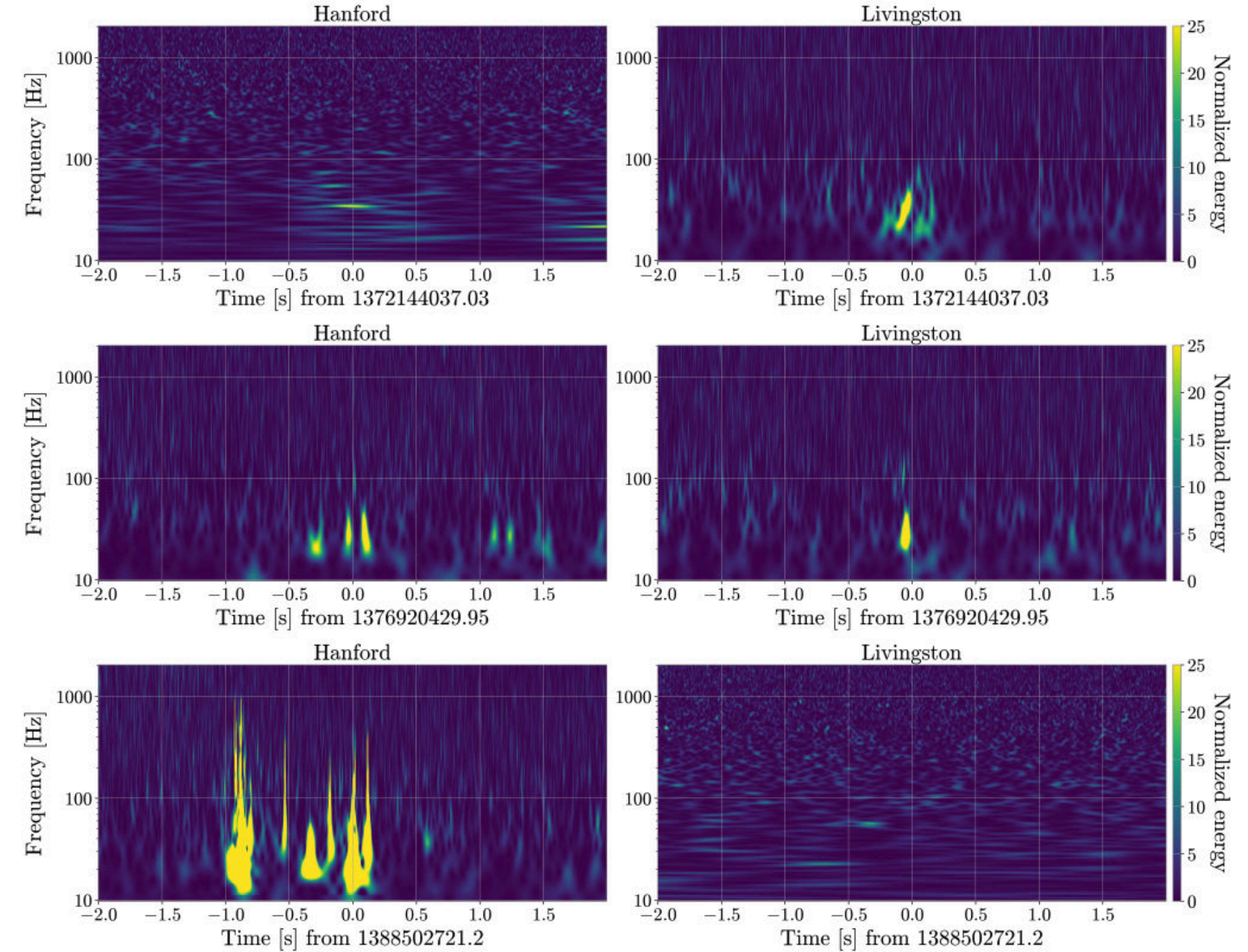
Version 4.0 of the Gravitational-Wave Transient Catalog (GWTC-4.0) adds new candidates detected by the LIGO, Virgo, and KAGRA observatories through the first part of the fourth observing run (O4a: 2023 May 24 15:00:00 to 2024 January 16 16:00:00 UTC) and a preceding engineering run. In these new data, we find 128



**Figure 2.** The marginal probability distributions for the source frame chirp mass  $\mathcal{M}$ , mass ratio  $q$ , effective inspiral spin  $\chi_{\text{eff}}$ , effective precession spin  $\chi_p$ , and luminosity distance  $D_L$  for O4a candidates with FAR  $< 1 \text{ yr}^{-1}$ . The colored upper half of the plot shows the marginal posterior distributions using our default agnostic priors (Abac et al. 2025b), while the white lower halves show these marginal distributions after reweighting according to the inferred population model (Abac et al. 2025g) for each BBH. The two NSBH candidates have only the non-reweighted posterior distributions. The vertical thickness of each region is proportional to the marginal posterior probability at that value for each candidate.

GW inference is fundamentally Bayesian:

$$p(\theta|d) \propto p(d|\theta)p(\theta)$$



**Figure 9.** Spectrograms for three candidates of interest. *Top panels:* Spectrograms for GW230630\_070659, a candidate identified by GSTLAL with FAR  $< 1 \text{ yr}^{-1}$  but which event validation indicates is likely of instrumental origin. We do not carry out parameter estimation on this candidate. *Middle panels:* Spectrograms for GW230824\_135331, a cWB-BBH candidate which is not identified as a candidate or subthreshold candidate by any matched filter search. Excess power is visible in LHO around the time of the candidate. *Bottom panels:* Spectrograms for GW240105\_151143, a candidate identified by only PYCBC. Significant excess power is present in LHO at the time of the candidate.



Astrophysical  
System

Understanding...



Science  
Investigations

Level 3 data

Level 2 data

Sources and  
Catalogues

Fitting  
Models

Gravitational Wave  
Data

Global Fit

Level 1 data

Doppler Shift  
Measurements

Processing  
Spacecraft Signals

Level 0 data

# LISA

## MISSION CONCEPT

# Scientific Context

- **The Next Generation (2030s):**
  - Einstein Telescope (ET).
  - Laser Interferometer Space Antenna (LISA).
- **The Leap:** Order of magnitude more sensitive.
- **The Problem:** The "Confusion Foreground".
  - Overlapping signals (thousands per year).
  - Non-stationary noise (Glitches).
  - Computational bottleneck for traditional MCMC.
- **MANGO Goal:** A robust ML framework for GW *separation and inference* using generative AI.



# Consortium

## CAD Team

**A. Rasamoela:** ML/Signal Processing  
**S. Caillou:** ML/Parameter Estimation  
**M. Dubois:** ML/Datasets  
**M. Pigou:** Infrastructure

## GW Team

**S. Marsat:** Datasets/Benchmarking  
**N. Tamanini:** Cosmology  
+ Postdocs

## Expected Funding

- ANR AAPG 2026
- 1 PhD Student
- 4 FTE Engineer/PostDoc

## Expected Impact

- First unified ML pipeline for ET/LISA.
- Solving the "Confusion" problem via Diffusion Model.
- Fast, calibrated posteriors via Latent Conditional Flow Matching.

## WP2 (Signals)

Objective: Denoising, Glitch Mitigation and Separation.

- Instead of just removing noise, we train ML models to learn the *manifold* of overlapping astrophysical signals.
- **Role 1 (Denoising):** Map corrupted strain  $\rightarrow$  clean signals (posterior samples).
- **Role 2 (Feature Extraction):** The internal representation (Latent  $z$ ) becomes the input for WP3.

## WP1 (Datasets)

- Creation of "Digital Twins" for ET & LISA.
- Integration with LDC (LISA Data Challenges) and MDC.
- Provides the massive training sets for WP2/WP3.

## WP4 (Populations)

- Scaling up from single events to populations.
- Likelihood-free inference of mass/spin distributions.
- Detecting outliers (New Physics).

## WP3 (Parameters)

Objective: Amortized inference for overlapping signals.

- **Method: Conditional Flow Matching (CFM)**
  - Traditional Normalizing Flows are unstable.
  - We use *Optimal Transport* paths conditioned on the features from WP2.
- **Input:** Latent  $z$  from WP2 (Denoising Model).
- **Process:** Regress a straight vector field from Prior to Posterior.
- **Advantage:** Simulation-free training, faster than MCMC.

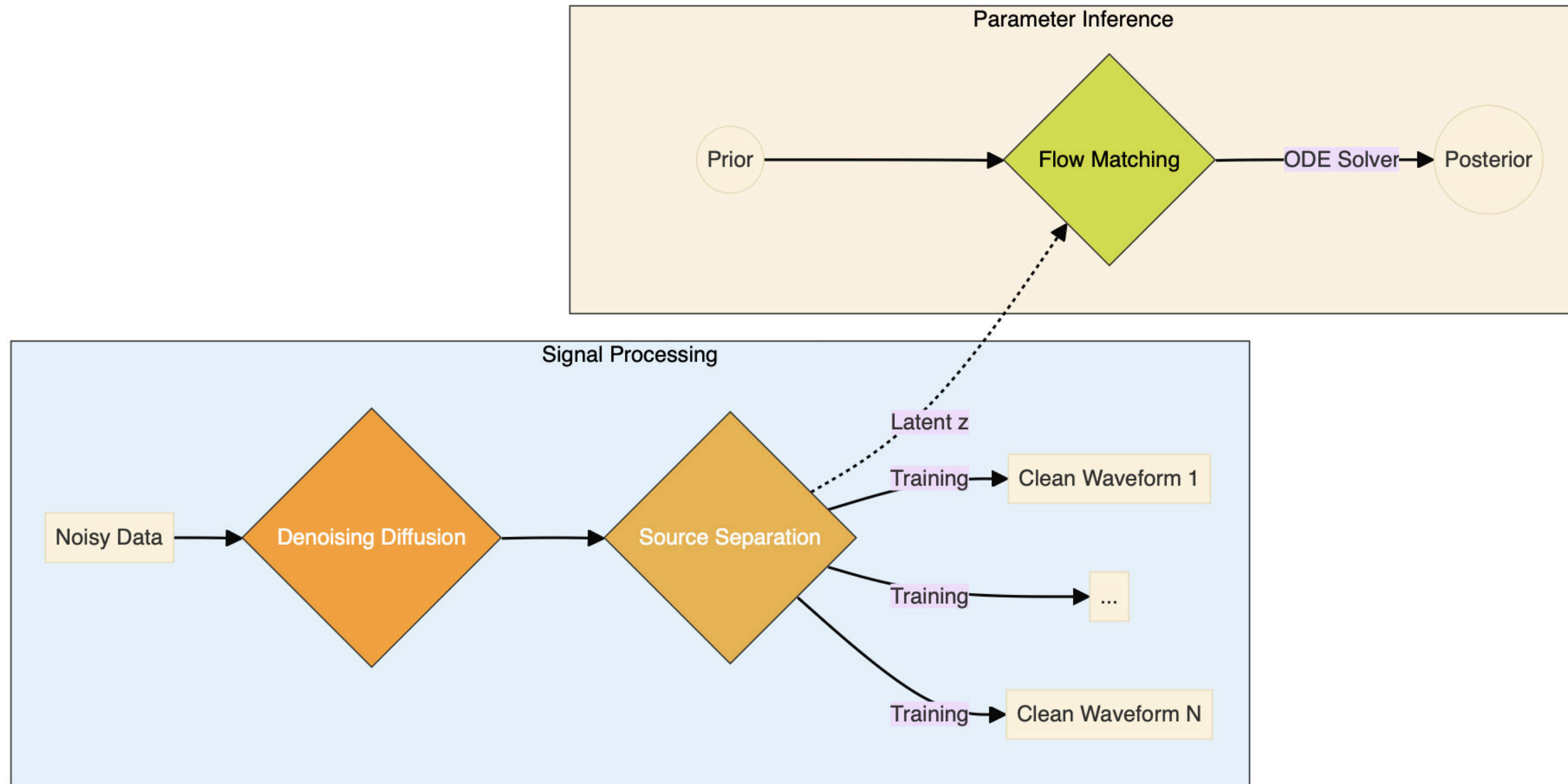
## WP5 (Infrastructure)

From Research to Production

- **Pipeline Integration:** Making ML models deployable in official pipelines (LISA DDPC).
- **HPC Deployment:** Training on national supercomputers (CC-IN2P3 & Jean-Zay).
- **Hardware:** Project includes funding for dedicated GPU Servers (4x H100 class) to handle the generative model training loads.



# Proposed Hybrid Architecture



# Denoising Diffusion Restoration Model

## Algorithm: Pre-training

**Input:** Dataset  $\mathcal{D}$  (Clean Waveforms + Glitches)

1. Sample  $d \sim \mathcal{D}$ ,  $t \sim \mathcal{U}(0, 1)$ ,  $\epsilon \sim \mathcal{N}(0, I)$
2. Corrupt:  $d_t = \sqrt{\bar{\alpha}_t}d + \sqrt{1 - \bar{\alpha}_t}\epsilon$
3. Predict Noise:  $\hat{\epsilon} = \epsilon_\phi(d_t, t)$
4. Optimization:  $\min_\phi ||\epsilon - \hat{\epsilon}||^2$

**Output:** Robust feature extractor  $\epsilon_\phi$  capable of filtering glitches.



# Flow Matching

## Algorithm: Conditional Flow Matching

**Frozen Input:** SCNet latent code  $z = E_{\text{WP2}}(d)$

1. Sample Prior  $\theta_0 \sim \mathcal{N}(0, I)$
2. Sample Target  $\theta_1 \sim p(\theta_{\text{true}})$
3. Interpolate:  $\theta_t = (1 - t)\theta_0 + t\theta_1$
4. Loss:  $\|v_\psi(\theta_t, t, z) - (\theta_1 - \theta_0)\|^2$

**Inference:** Solve ODE  $d\theta/dt = v_\psi$  starting from noise.



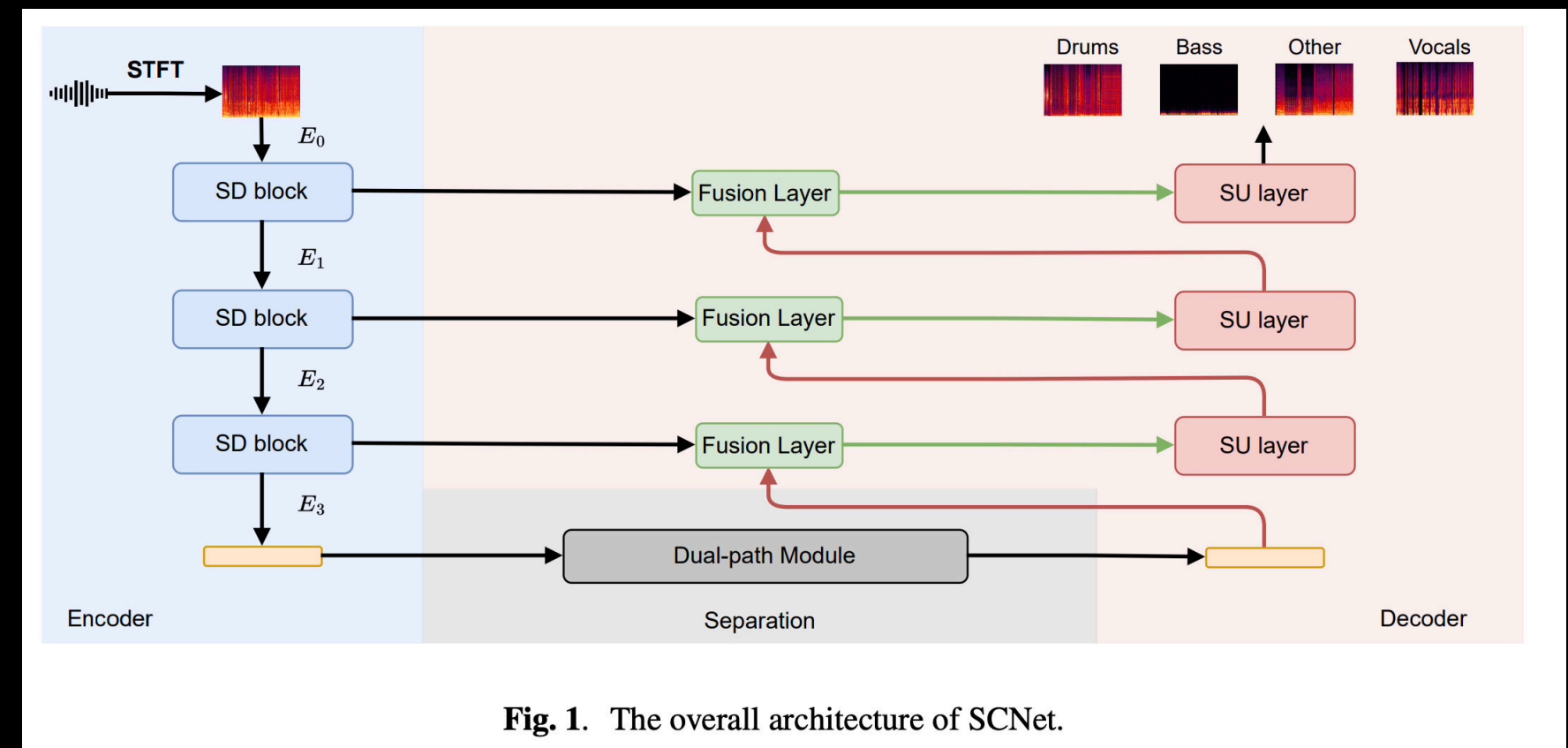
# GWINESS

GRAVITATIONAL WAVE INFERENCE  
NEURAL SOURCE SEPARATION





- **GWINESS** project started in January 2025
- Blind source separation of overlapping GWs
- Inspired by **music/speech separation**
- **SCNet** deep learning architecture for music source separation in Time-Freq Domain



**Fig. 1.** The overall architecture of SCNet.

# Motivation

- LISA will detect **many overlapping GW** signals from different type of source
- Classical Bayesian inference methods are **computationally expensive**
- Explore deep learning for efficient **source separation and sampling**

## The global fit

Scores from a penguin cacophony



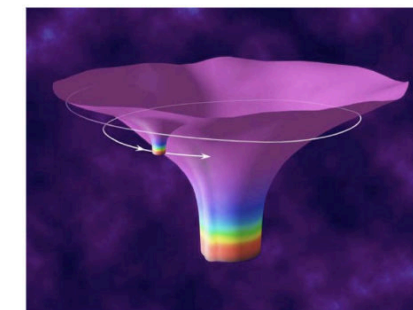
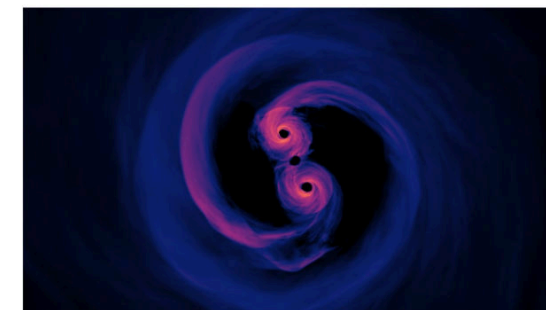
Hundreds



Tens to thousands



Millions

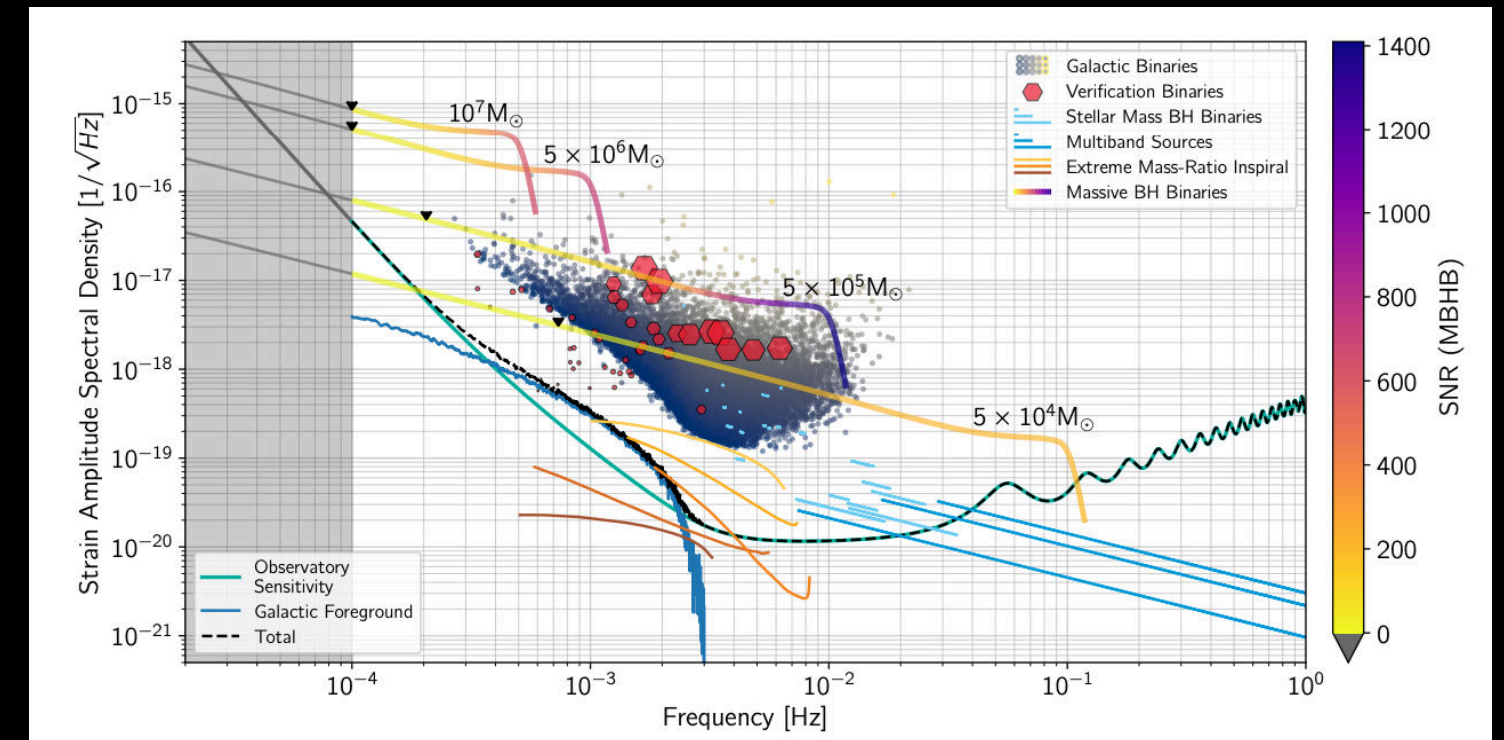


from R. Buscicchio's talk, Toulouse, 10/2024



# Challenges in LISA

- Complex noise structures (+glitches and gaps)
- Large number of overlapping sources
- High dimensionality of the model
- Presence of correlated parameters



<https://arxiv.org/pdf/2402.07571>

- GBs : narrow band signals in the frequency domain
- BHBs: transient signals (response varying both with time and freq)
- EMRIs: complex long-lasting signals with low SNR and timescales of several years



Data Input

Drop HDF5 files here

or

Browse Files

Loaded File:

Name: mixed\_test.hdf5

Size: 6.43 GB

Modified: 29/04/2025 16:35:38

Analysis Tools

Spectrogram

Parameter

Correlation

Source Information

MBHB

GW150914-like

2.3 Gpc

EMRI

Extreme Mass Ratio

4.1 Gpc

Galactic

White Dwarf Binary

8 kpc

Visualization Controls

Item: 10

Start

Stop

Fullscreen

Mixed Input

Export

MBHB

Focus

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

Sampling Rate: 0.1 Hz

Separated Outputs

Hide

All

Show

MBHB

Massive Black Hole Binary

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

SNR: 24.5

Frequency: 0.01-0.1 Hz

EMRI

Extreme Mass Ratio Inspiral

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

SNR: 18.2

Frequency: 0.001-0.01 Hz

Galactic

White Dwarf Binary

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

SNR: 32.7

Frequency: 0.1-1 Hz

Residual

Remaining Noise

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

Power: 0.05

Frequency: Full Band

Data Input

Drop HDF5 files here

or

Browse Files

Loaded File:

Name: mixed\_test.hdf5

Size: 6.43 GB

Modified: 29/04/2025 16:35:38

Analysis Tools

Spectrogram

Parameter

Correlation

Source Information

MBHB

GW150914-like

2.3 Gpc

EMRI

Extreme Mass Ratio

4.1 Gpc

Galactic

White Dwarf Binary

8 kpc

Visualization Controls

Item: 42

Start

Stop

Fullscreen

Mixed Input

Export

MBHB

Focus

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

Sampling Rate: 0.1 Hz

Separated Outputs

Hide

All

Show

MBHB

Massive Black Hole Binary

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

SNR: 24.5

Frequency: 0.01-0.1 Hz

EMRI

Extreme Mass Ratio Inspiral

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

SNR: 18.2

Frequency: 0.001-0.01 Hz

Galactic

White Dwarf Binary

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

SNR: 32.7

Frequency: 0.1-1 Hz

Residual

Remaining Noise

Time Domain

Frequency (Hz)

Time (days)

Log Amplitude

Power: 0.05

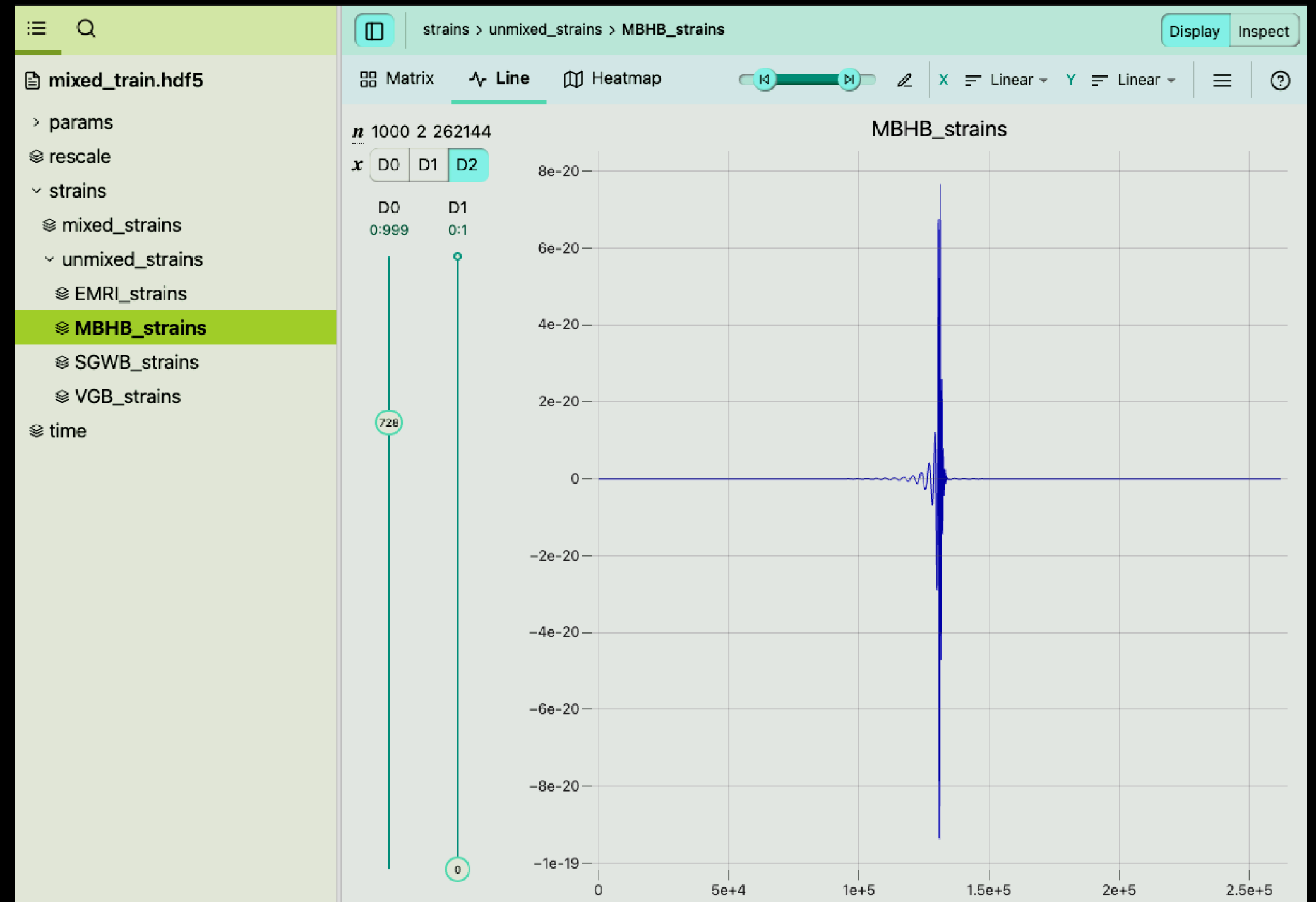
Frequency: Full Band

14



# Dataset generation

**generate - sample - mix**  
**1000 mixed samples**  
**16 GB (disk)**



```
(denoise_few) [arasamoela@ccwgislurm0200 demos]$ python generate_mixture.py 1000 datasets/training
Generating VGB: 100%|██████████████████████████████████████████████████████████████████████████| 43/43 [00:01<00:00, 40.41it/s]
Generating SGWB: 100%|██████████████████████████████████████████████████████████████████████████| 10000/10000 [00:19<00:00, 509.62it/s]
Generating MBHB: 100%|██████████████████████████████████████████████████████████████████████████| 1000/1000 [03:47<00:00, 4.40it/s]
Generating mixed dataset with 100 samples, TDI channel: AE, timestamp: 50961046
Generating mixed datasets: 100%|██████████████████████████████████████████████████████████████████████████| 100/100 [00:00<00:00, 313.64it/s]
Writing to file ----> datasets/training/50961046/mixed_test.hdf5
Generating mixed dataset with 1000 samples, TDI channel: AE, timestamp: 50961046
Generating mixed datasets: 100%|██████████████████████████████████████████████████████████████████████████| 1000/1000 [00:03<00:00, 324.48it/s]
Writing to file ----> datasets/training/50961046/mixed_train.hdf5
(denoise_few) [arasamoela@ccwgislurm0200 demos]$ ^C
```

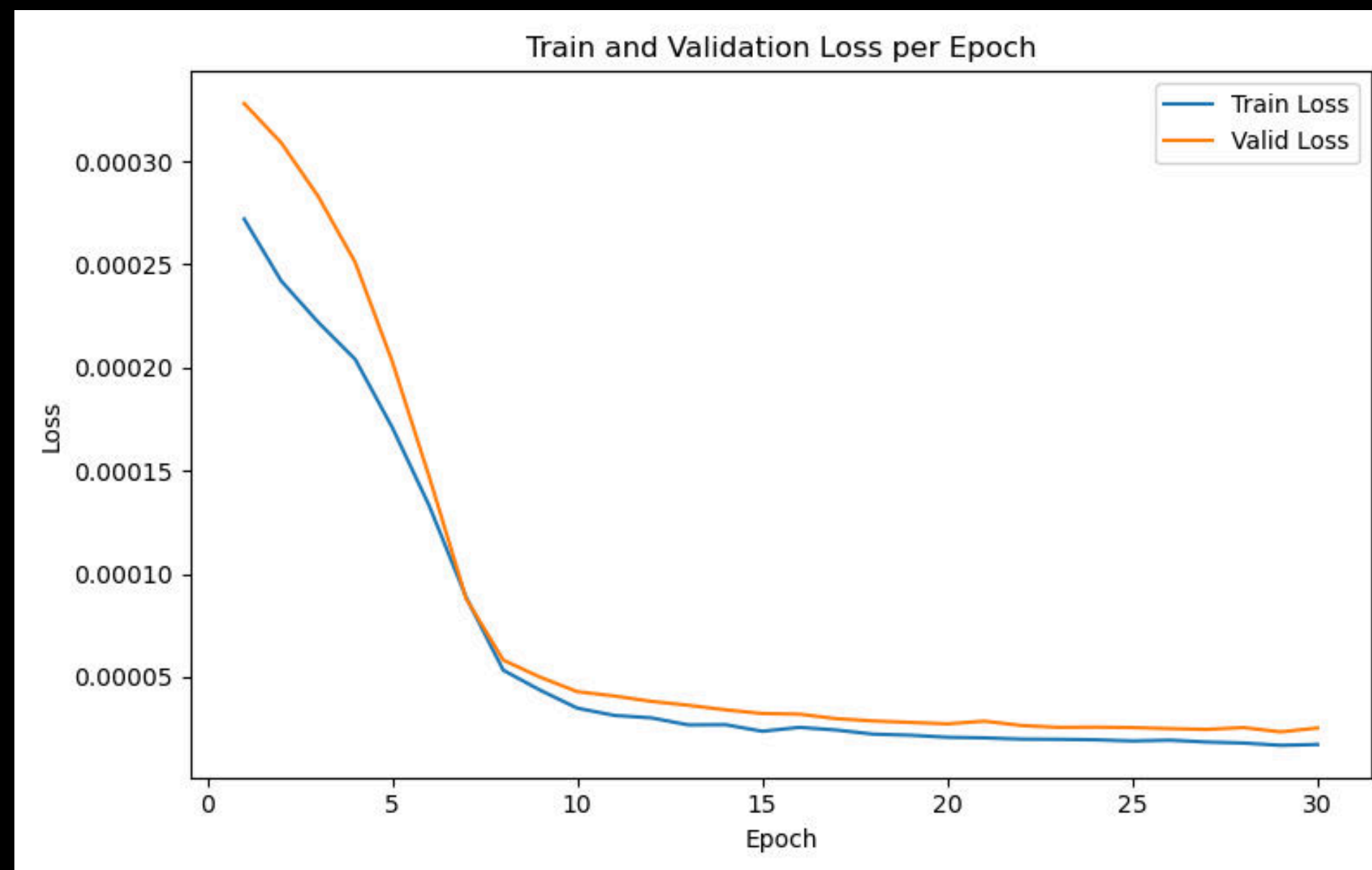


# Training

50 epoch -> 1 day  
42M parameters  
65 GB (GPU)

```
025-04-16 10:36:07,577 - INFO - Learning rate adjusted to 0.0003
025-04-16 10:36:07,579 - INFO - -----
025-04-16 10:36:07,579 - INFO - Training Epoch 1 ...
025-04-16 10:45:18,159 - INFO - Train Summary | Epoch 1 | Loss=0.0003 | Grad=6.9445
025-04-16 10:45:18,160 - INFO - -----
025-04-16 10:45:18,160 - INFO - Cross validation...
025-04-16 10:51:40,228 - INFO - Valid Summary | Epoch 1 | Loss=0.0003 | Nsdr=-16.867
025-04-16 10:51:40,986 - INFO - Learning rate adjusted to 0.0003
```

```
2025-04-17 13:24:13,898 - INFO - Cross validation...
2025-04-17 13:28:28,931 - INFO - Valid Summary | Epoch 52 | Loss=0.0000 | Nsdr=-4.503
2025-04-17 13:31:55,577 - INFO - Total number of parameters: 42181232
2025-04-17 13:36:34,662 - INFO - train/valid set size: 1000 500
2025-04-17 13:36:34,729 - INFO - Loading checkpoint model: datasets/save/checkpoint.th
2025-04-17 13:53:37,565 - INFO - Total number of parameters: 42181232
2025-04-17 13:55:35,552 - INFO - Total number of parameters: 42181232
2025-04-17 14:00:21,909 - INFO - train/valid set size: 1000 500
2025-04-17 14:00:25,497 - INFO - Learning rate adjusted to 0.0003
```



## 2.5. Loss function

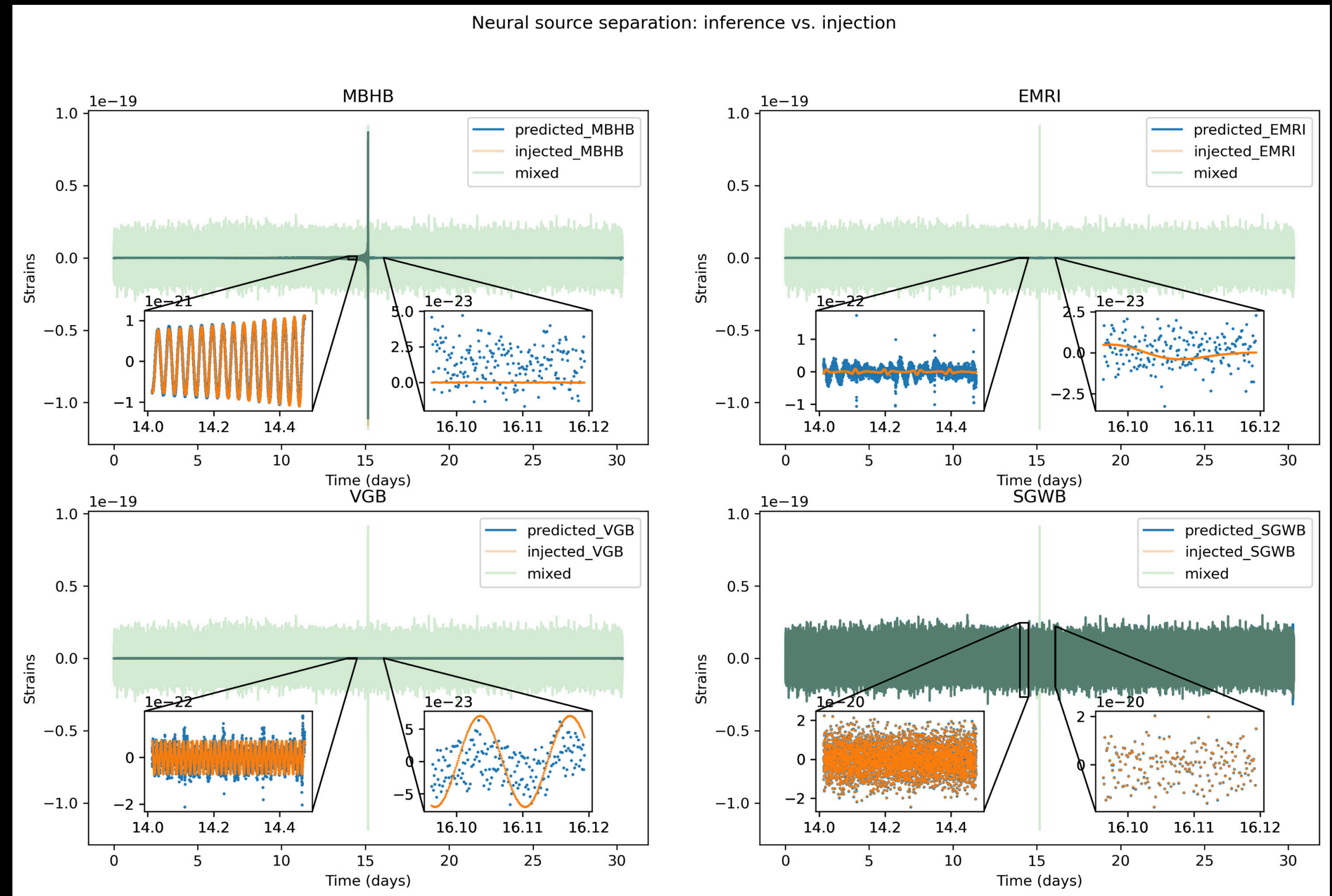
Previous work generally considers waveform similarity as the training target. But this is not so closely related to the spectrogram. So we use the root mean squared error (RMSE) loss of complex-valued spectrogram as the loss function:

$$\mathcal{L}_s = \sqrt{(r - \hat{r})^2 + (i - \hat{i})^2} \quad (4)$$

where  $\hat{r}$  and  $\hat{i}$  represent the real and imaginary parts of the source spectrogram respectively.

# Inference

1 year  $\rightarrow$  2s  
1.1 GB (checkpoint in disk)  
3 GB (GPU)



---

# Future directions

While the **benefits** are clear, there are also many **challenges** to consider:

- working on a proof of concept
- !!! looking for collaboration !!!
- product assurance / acceptable AI

The **GWINESS** approach requires further investigation:

- dataset generation during training (many overlapping sources)
- hyperparameter tuning
- test with parameter estimation

**Data quality** and **computational resources** impact its effectiveness:

- implement fast waveform generator
- develop training dataset pipeline
- need a GPU-based cluster suited for large model training





# Thanks!

---

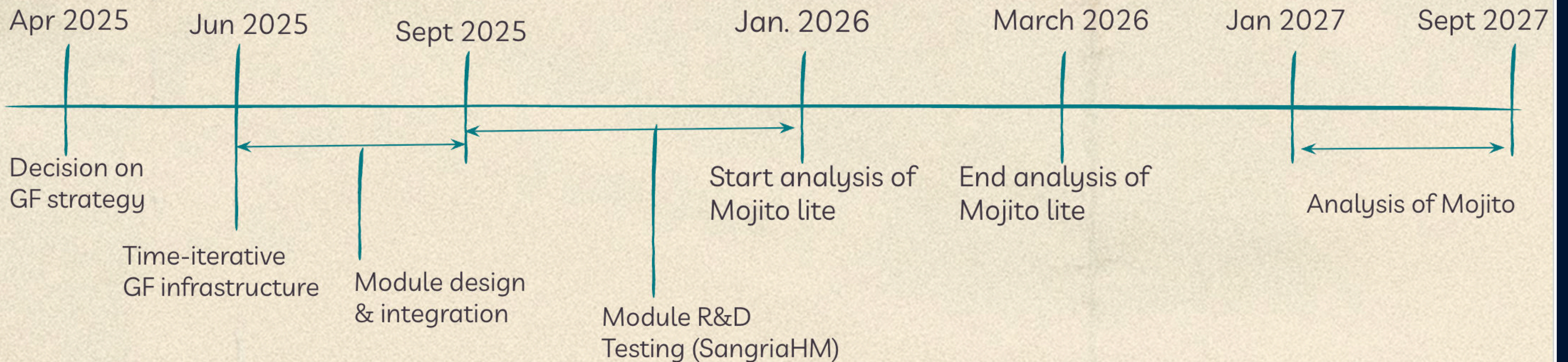
Do you have any questions?





# Integration to LISA roadmap

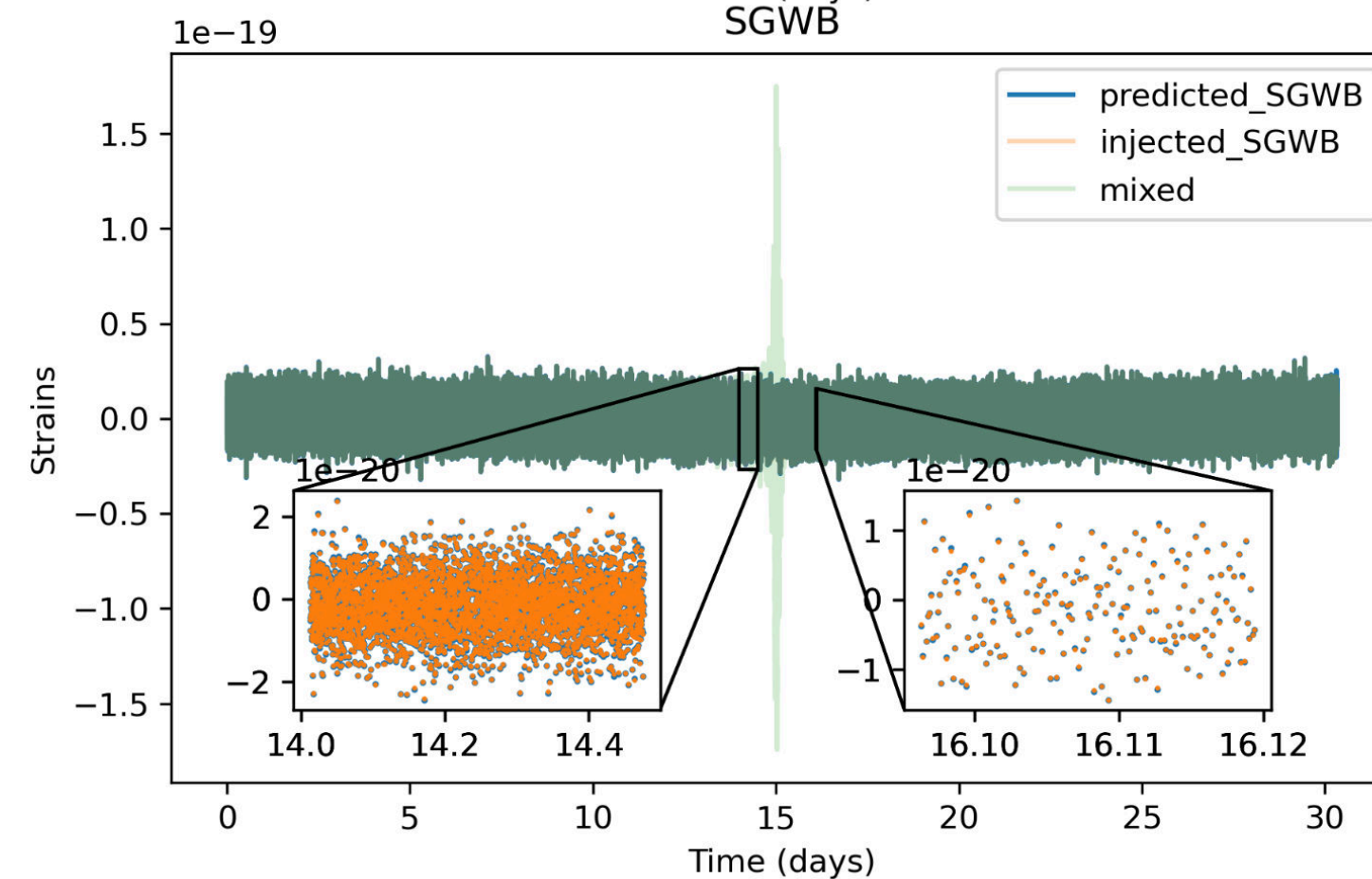
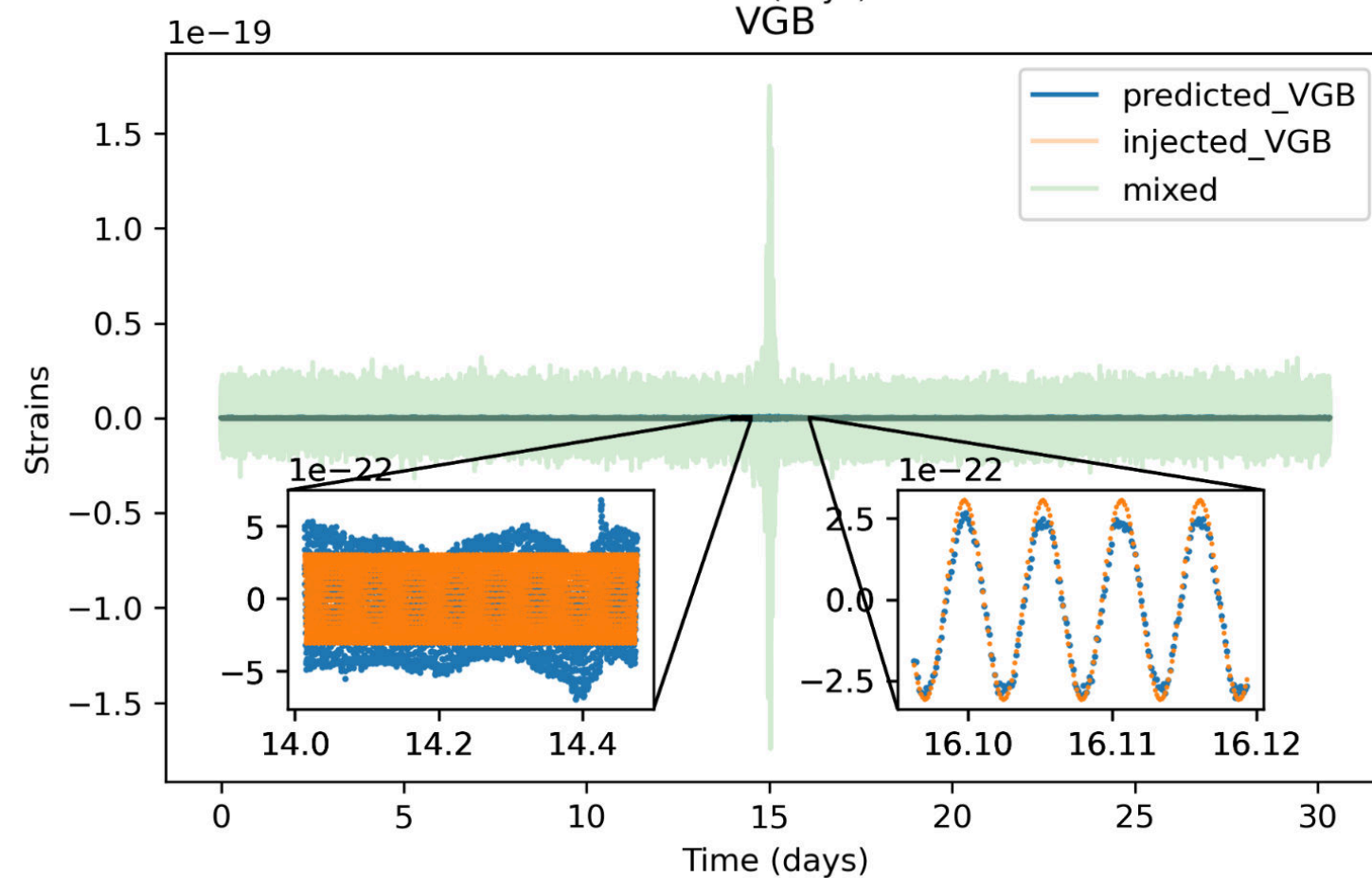
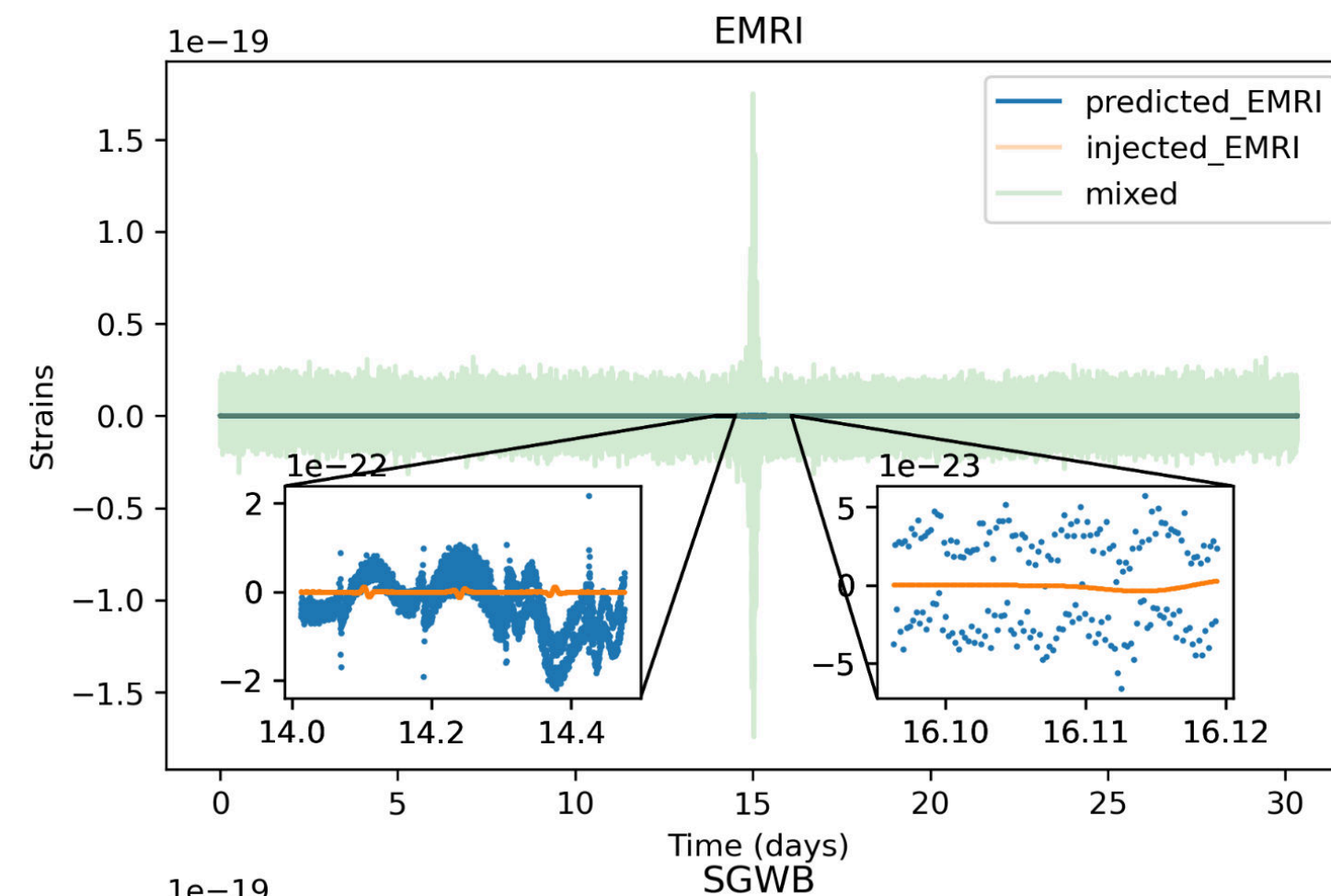
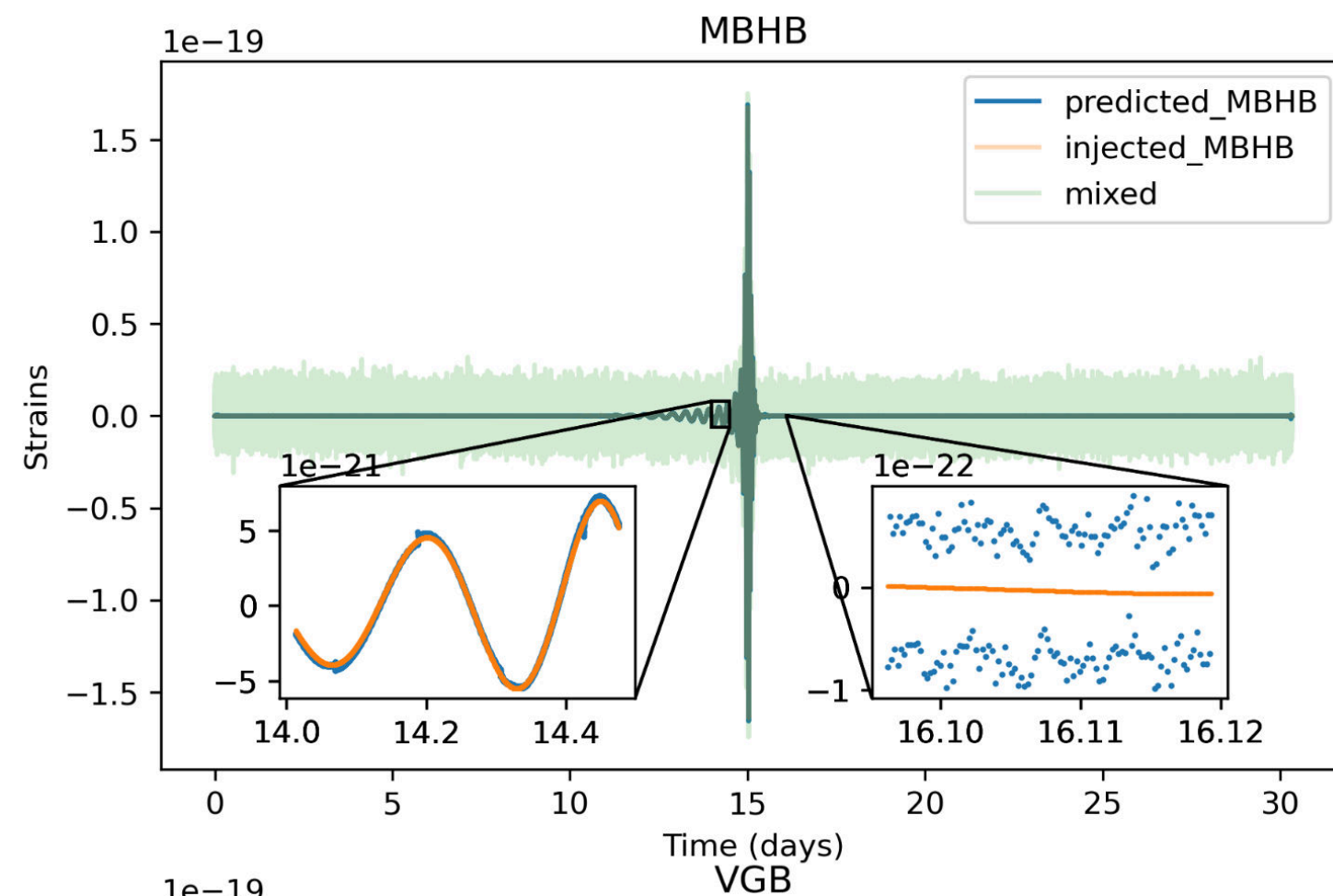
## Timeline





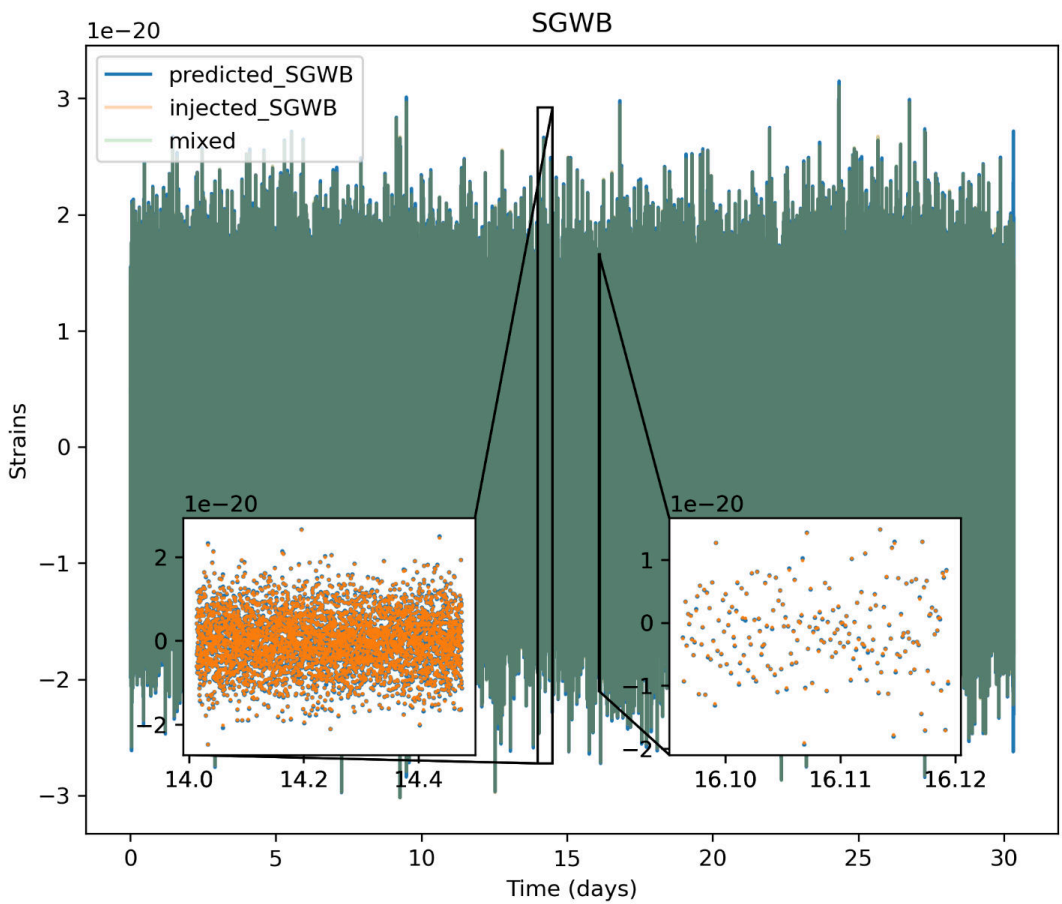
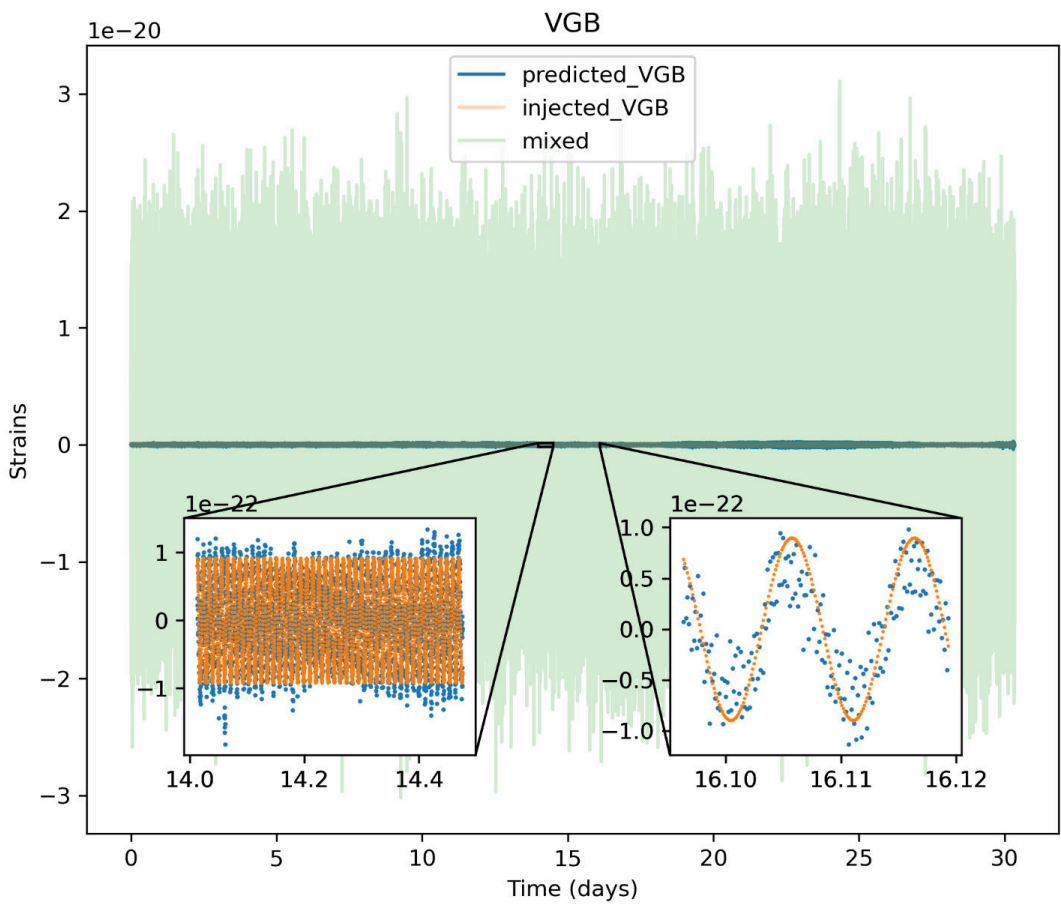
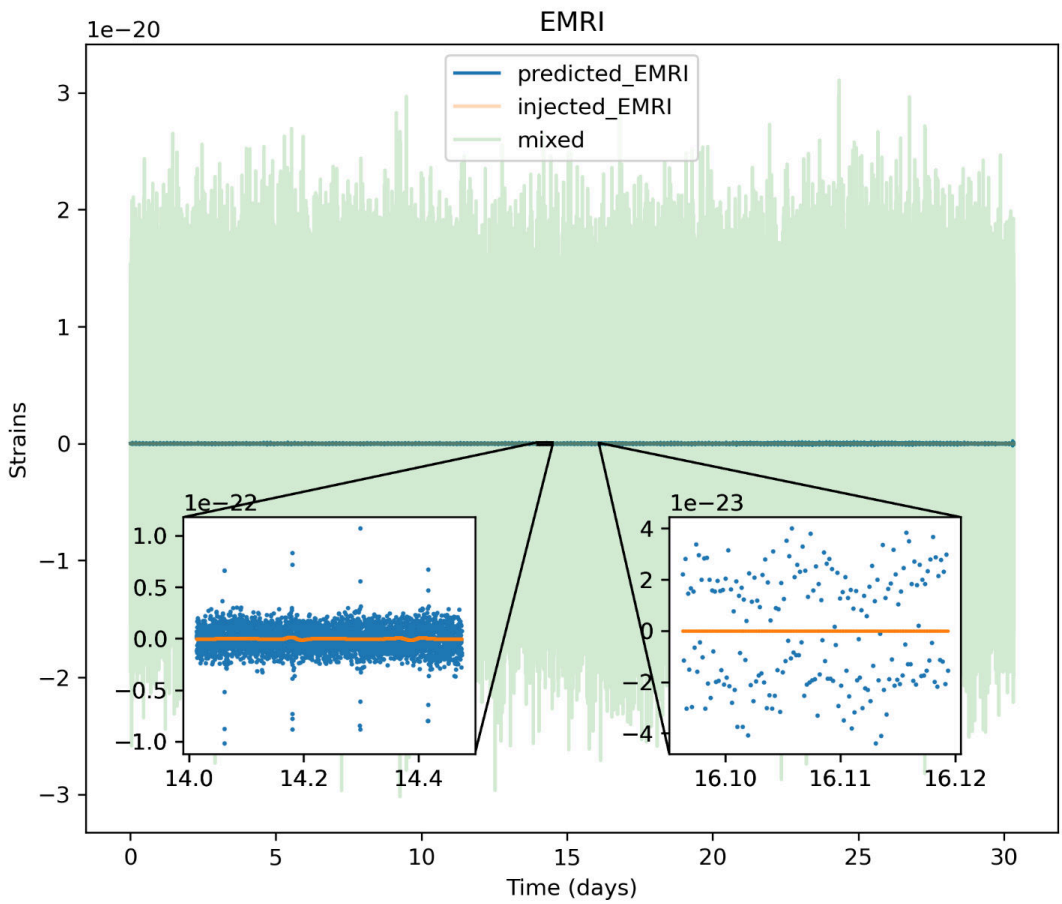
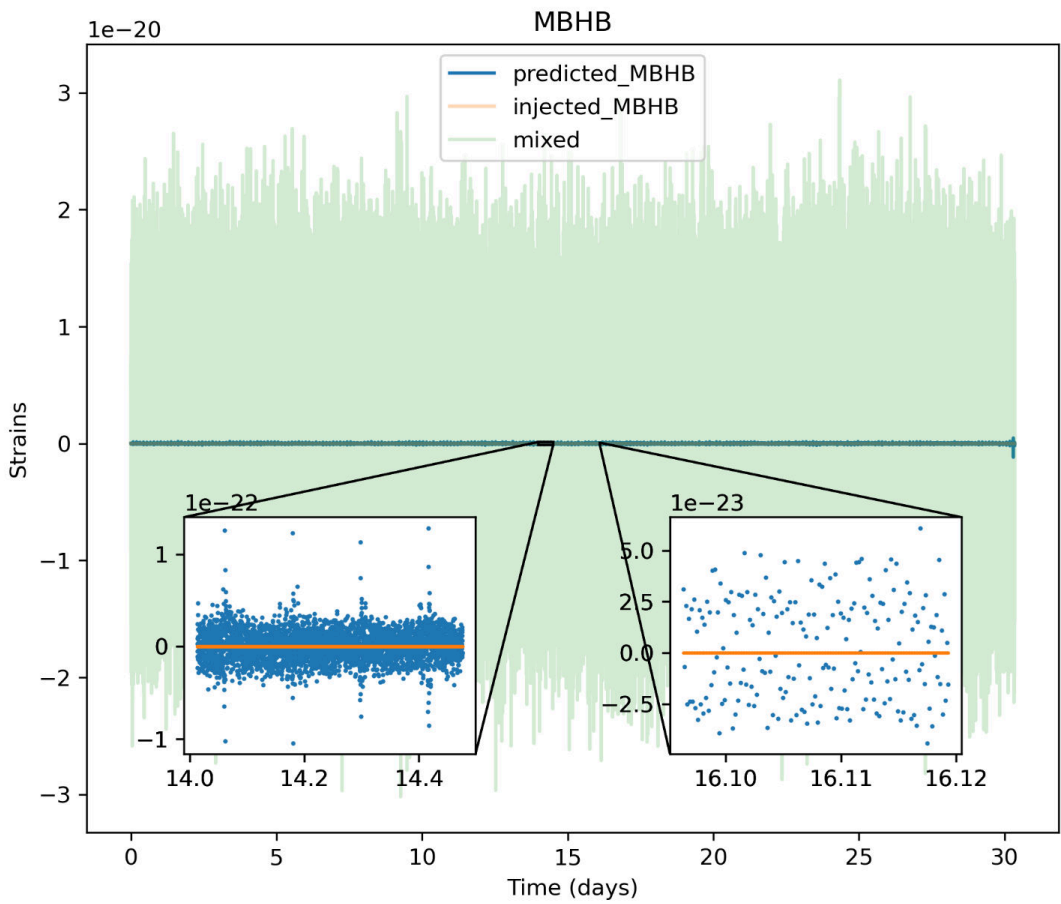
Neural source separation: inference vs. injection

With MBHB





# Without MBHB



# 1. Large class-level separation

LISA TDI → SCNet Coarse Separation → Source Type Channels

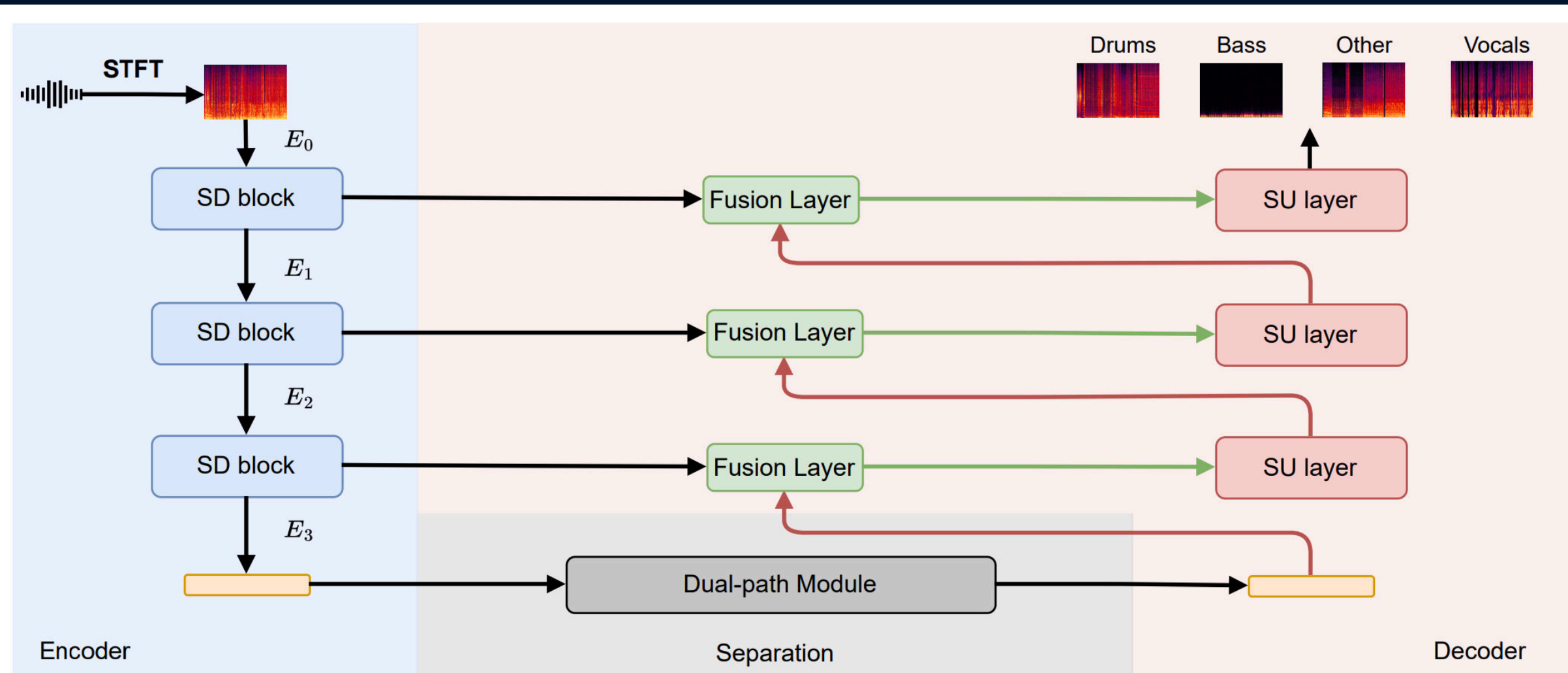


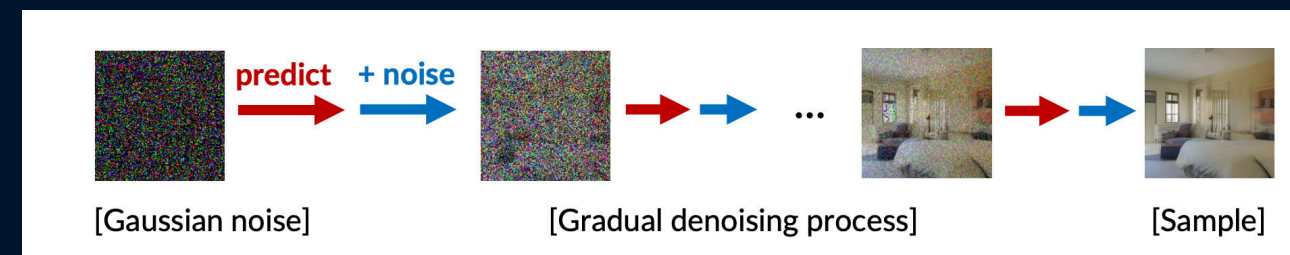
Fig. 1. The overall architecture of SCNet.

```
gwiness > coalescnet > conf > ! config.yaml
1 data:
2   mixed_train_path: /sps/l2it/arasamoela/gwiness/coalescnet/datasets/mixed_train.hdf5
3   mixed_valid_path: /sps/l2it/arasamoela/gwiness/coalescnet/datasets/mixed_test.hdf5
4   samplerate: 0.1
5   segment: 1
6   shift: 0.5
7   channels: 2
8   normalize: true
9   metadata:
10    sources: ['EMRI', 'MBHB', 'SGWB', 'VGB']
11
12 ema:
13   epoch: [0.9, 0.95]
14   batch: [0.9995, 0.9999]
15
16 model:
17   sources: ['EMRI', 'MBHB', 'SGWB', 'VGB']
18   audio_channels: 2
19   # Main structure
20   dims: [4, 64, 128, 256]
21   # STFT
22   nfft: 4096
23   hop_size: 1024
24   win_size: 4096
25   normalized: True
26   # SD/SU layer
27   band_SR: [0.225, 0.372, 0.403]
28   band_stride: [1, 4, 16]
29   band_kernel: [3, 4, 16]
30   # Convolution Module
31   conv_depths: [3, 2, 1]
32   compress: 4
33   conv_kernel: 3
34   # Dual-path RNN
35   num_dplayer: 6
36   expand: 1
37
38 epochs: 130
39 batch_size: 16
40
```



## 2. Probabilistic refinement

LISA TDI → SCNet Coarse Separation → Source Type Channels

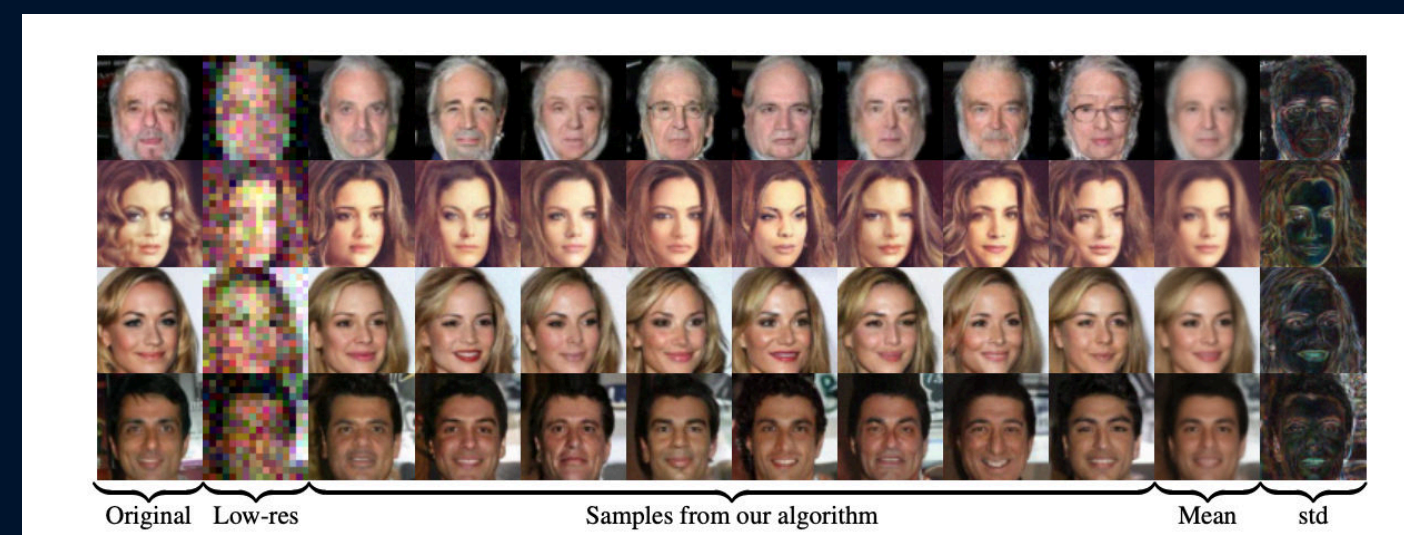


DDRM Denoising / Sampling



Clean Mixture Samples

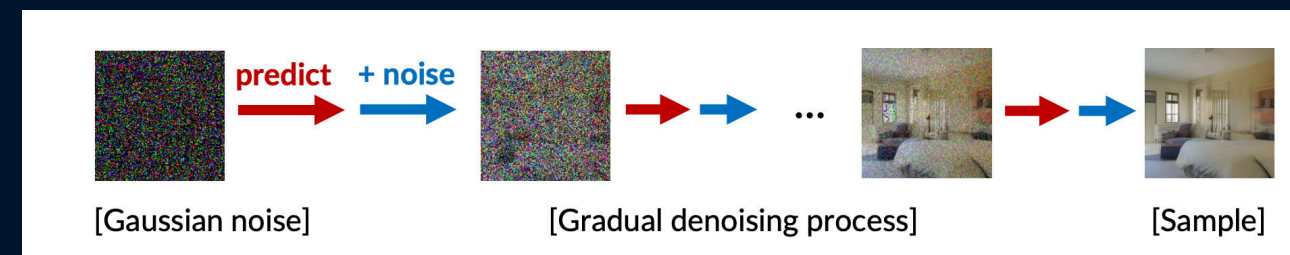
(per type, e.g., GBs within small frequency windows)





### 3. Overlapping sources unmixing

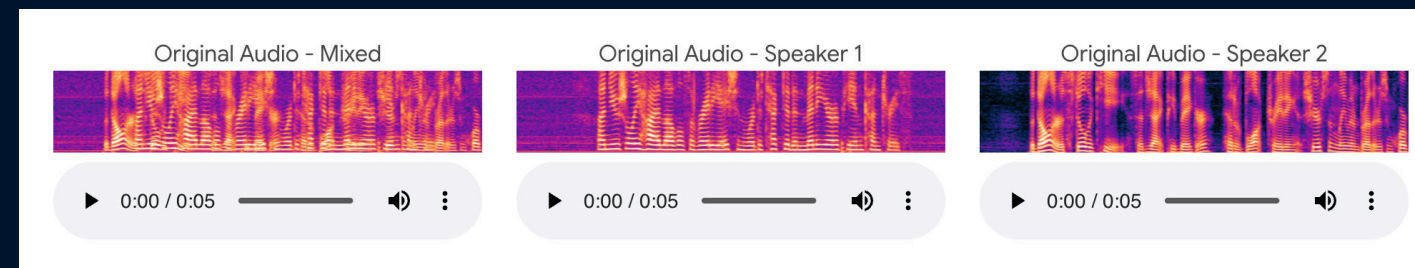
LISA TDI → SCNet Coarse Separation → Source Type Channels



DDRM Denoising / Sampling



Clean Mixture Samples



SepReformer Fine Separation



Clean Waveform Samples  $\{\text{source}_1, \dots, \text{source}_n\}$



# 3. Overlapping sources unmixing

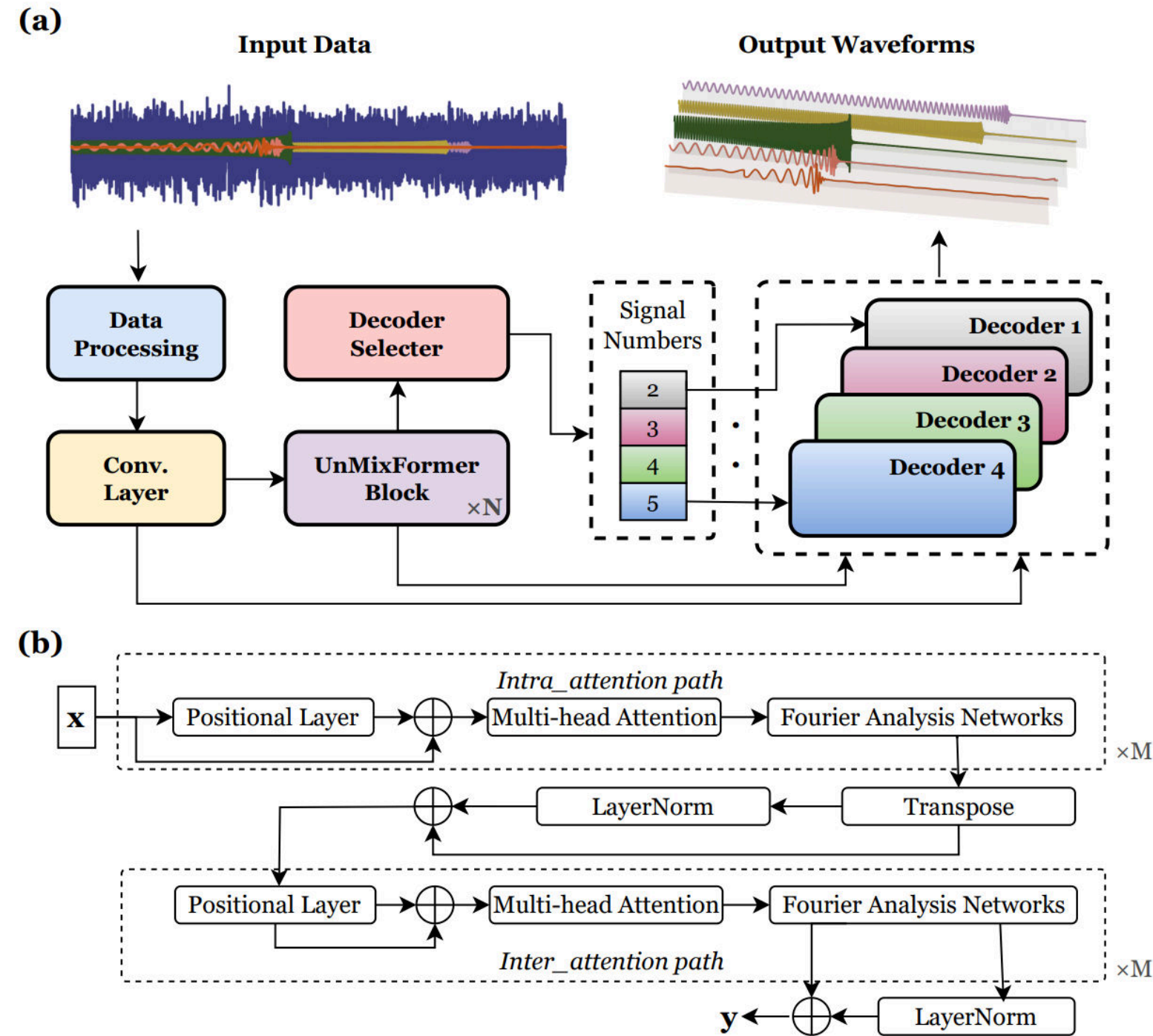


Figure 2. **UnMixFormer Architecture.** (a). The overall framework for counting and separating overlapping GW signals. We firstly employ CNN-based encoders to extract data embeddings, which are then fused and passed into UnMixFormer blocks. The counting head predicts the number of sources and activates the appropriate decoder to reconstruct individual waveforms. (b). The core UnMixFormer block operates with intra- and inter-attention mechanisms to capture fine-grained local features and global context. FAN layers in the feedforward module enhance periodic feature modeling and the positional encoding incorporates sequential information, enabling efficient separation of overlapping signals.

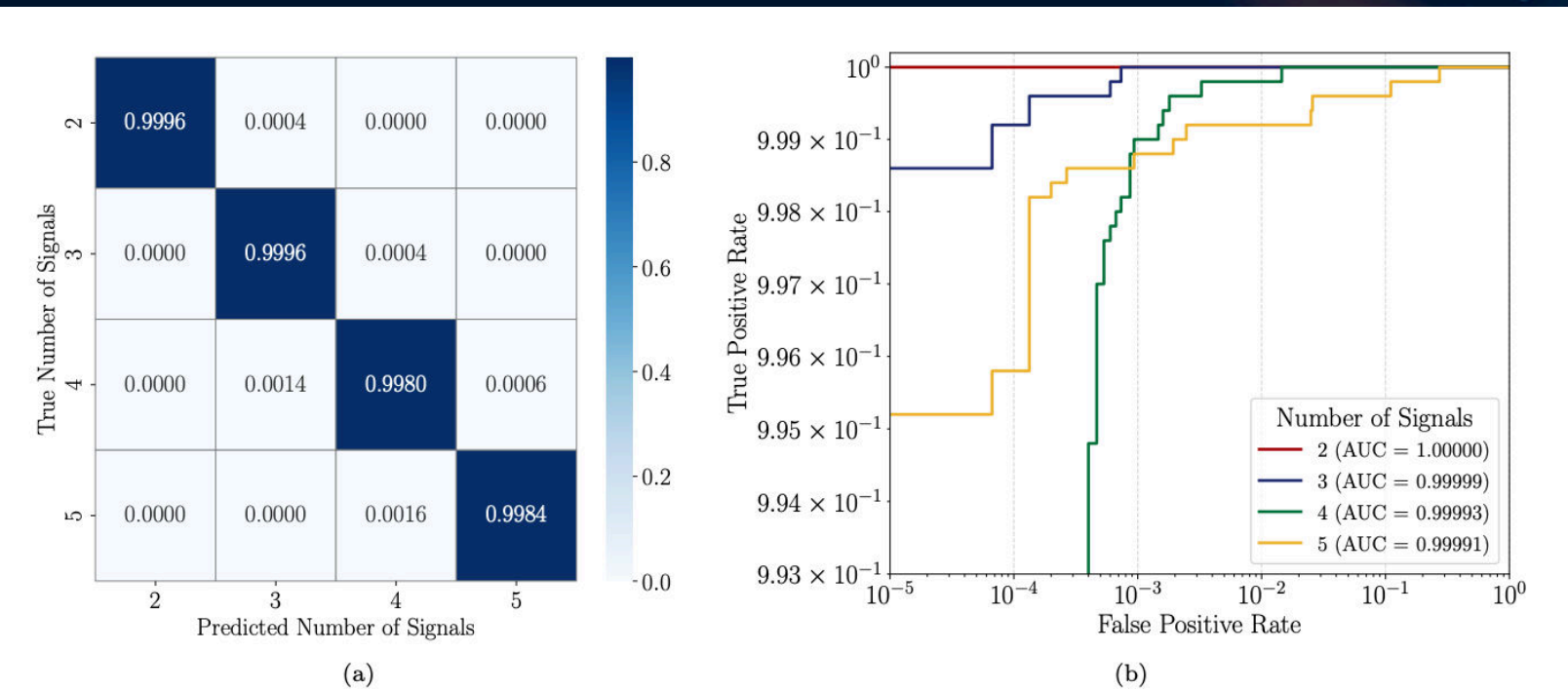


Figure 3. **Counting performance of overlapping CBC signals.** (a). The normalized confusion matrix shows the high accuracy of predicting the number of overlapping signals, with correct predictions dominating the diagonal entries (2 to 5 signals). (b). ROC curves illustrating the performance of signal counting for varying numbers of signals. The curves demonstrate near-perfect detection across all cases.

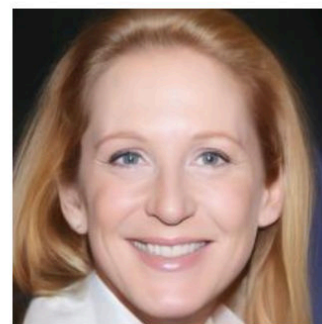
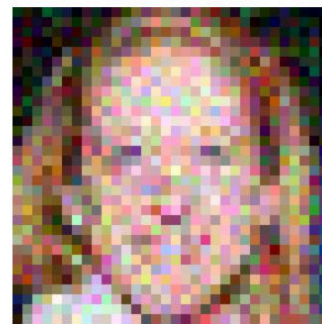
<https://arxiv.org/pdf/2412.18259>



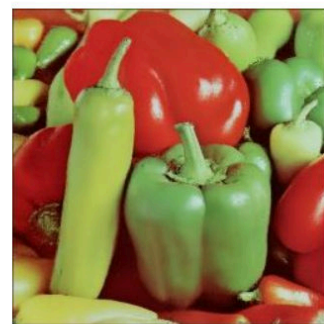
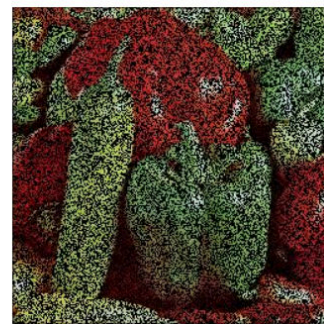
# Bayesian Deep Learning

## Denoising Diffusion Restoration Models

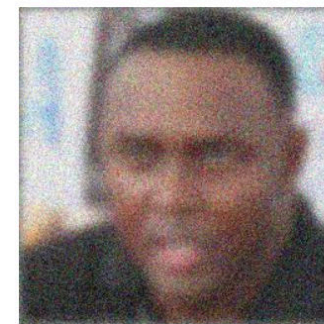
Super-resolution



Inpainting

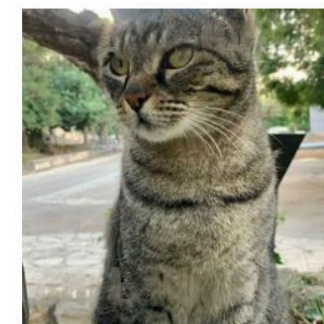


Deblurring



Observations  
(Inputs)  
 $y$

Outputs from  
our method  
 $x_0$



## Applicable to other domains as well!

### Speech

#### A VERSATILE DIFFUSION-BASED GENERATIVE REFINER FOR SPEECH ENHANCEMENT

Ryosuke Sawata Naoki Murata Yuhta Takida Toshimitsu Uesaka  
Takashi Shibuya Shusuke Takahashi Yuki Mitsufuji

Sony Group Corporation, Tokyo, Japan

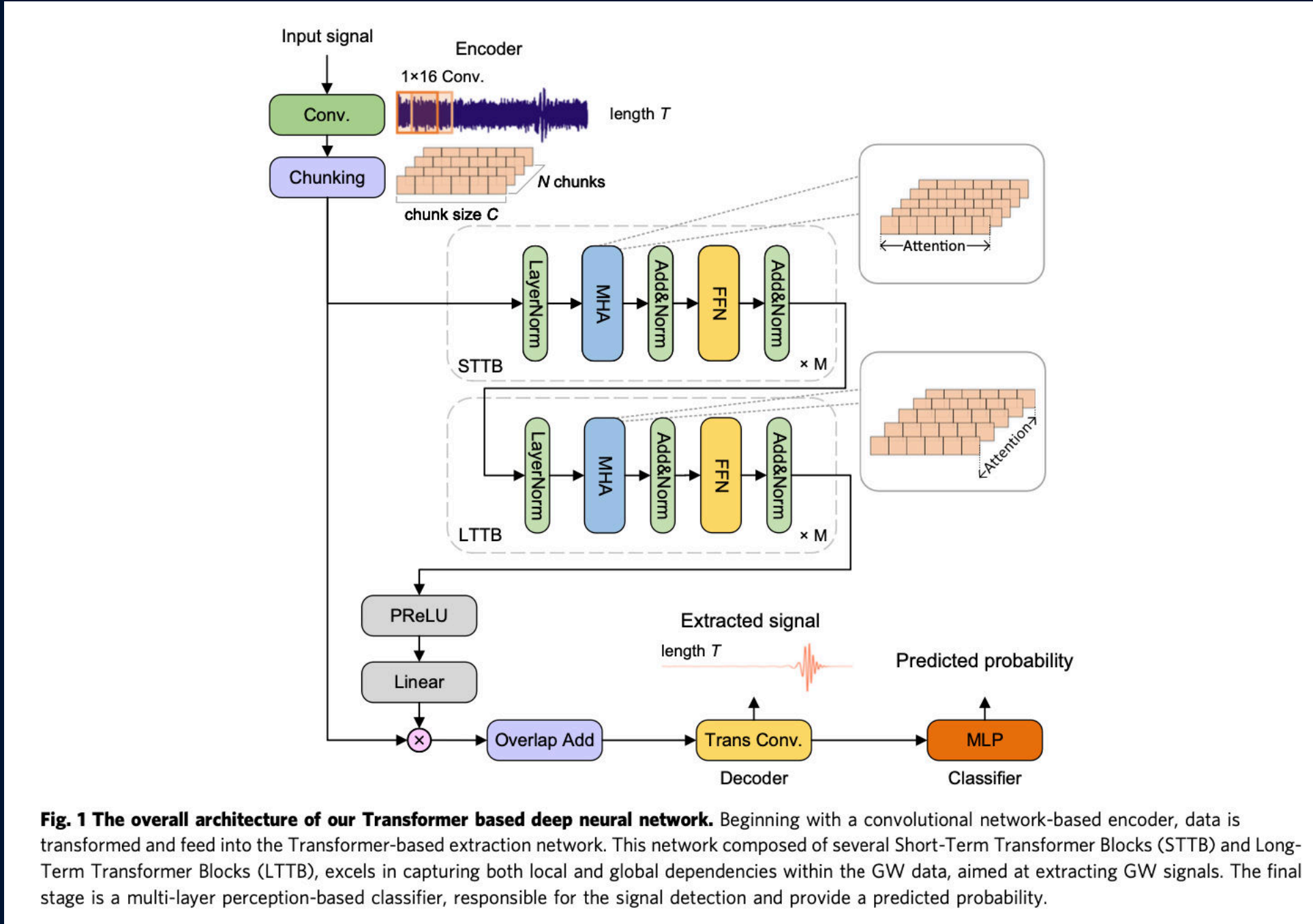
#### UNSUPERVISED VOCAL DEREVERBERATION WITH DIFFUSION-BASED GENERATIVE MODELS

Koichi Saito Naoki Murata Toshimitsu Uesaka Chieh-Hsin Lai  
Yuhta Takida Takao Fukui Yuki Mitsufuji

Sony Group Corporation, Tokyo, Japan



# Signal denoising



**Table 1 Summary of parameter setups in EMRI signal simulation.**

Parameter	Lower bound	Upper bound
$M$	$10^5 M_{\odot}$	$10^7 M_{\odot}$
$a$	$10^{-3}$	0.99
$e_0$	$10^{-3}$	0.5
$\cos i$	-1	1

**Table 2 Summary of parameter setups in MBHB signal simulation.**

Parameter	Lower bound	Upper bound
$M_{tot}$	$10^6 M_{\odot}$	$10^8 M_{\odot}$
$q$	0.01	1
$s_1^z$	-0.99	0.99
$s_2^z$	-0.99	0.99

**Table 3 Summary of parameter setups in BWD signal simulation.**

Parameter	$f$	$\dot{f}$
Range-1	[0.1, 4]mHz	$[-3 \times 10^{-17}, 6 \times 10^{-16}] \text{Hz}^2$
Range-2	[4, 15]mHz	$[-3 \times 10^{-15}, 4 \times 10^{-14}] \text{Hz}^2$



# Signal denoising & separation

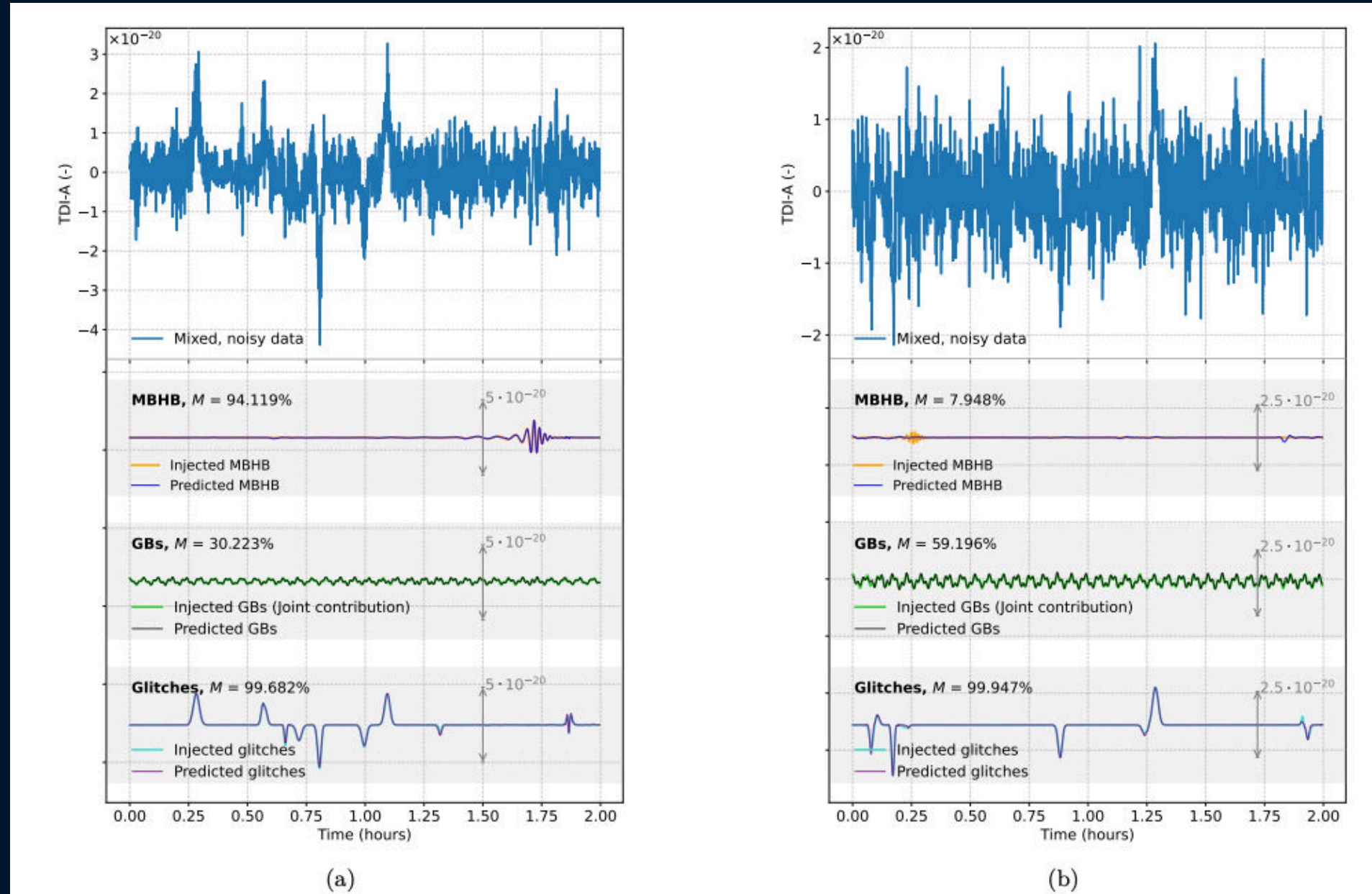


FIG. 7: Comparison of injected waveforms and model predictions for a low-amplitude MBHB merger buried in stationary noise. In panel (a), the deep source separation framework successfully detects and reconstructs the MBHB signal. However, in (b), where the signal amplitude is further diminished, the model fails. In such cases, further investigation is needed to determine whether the issue lies in the shared encoder or the MBHB decoder head.

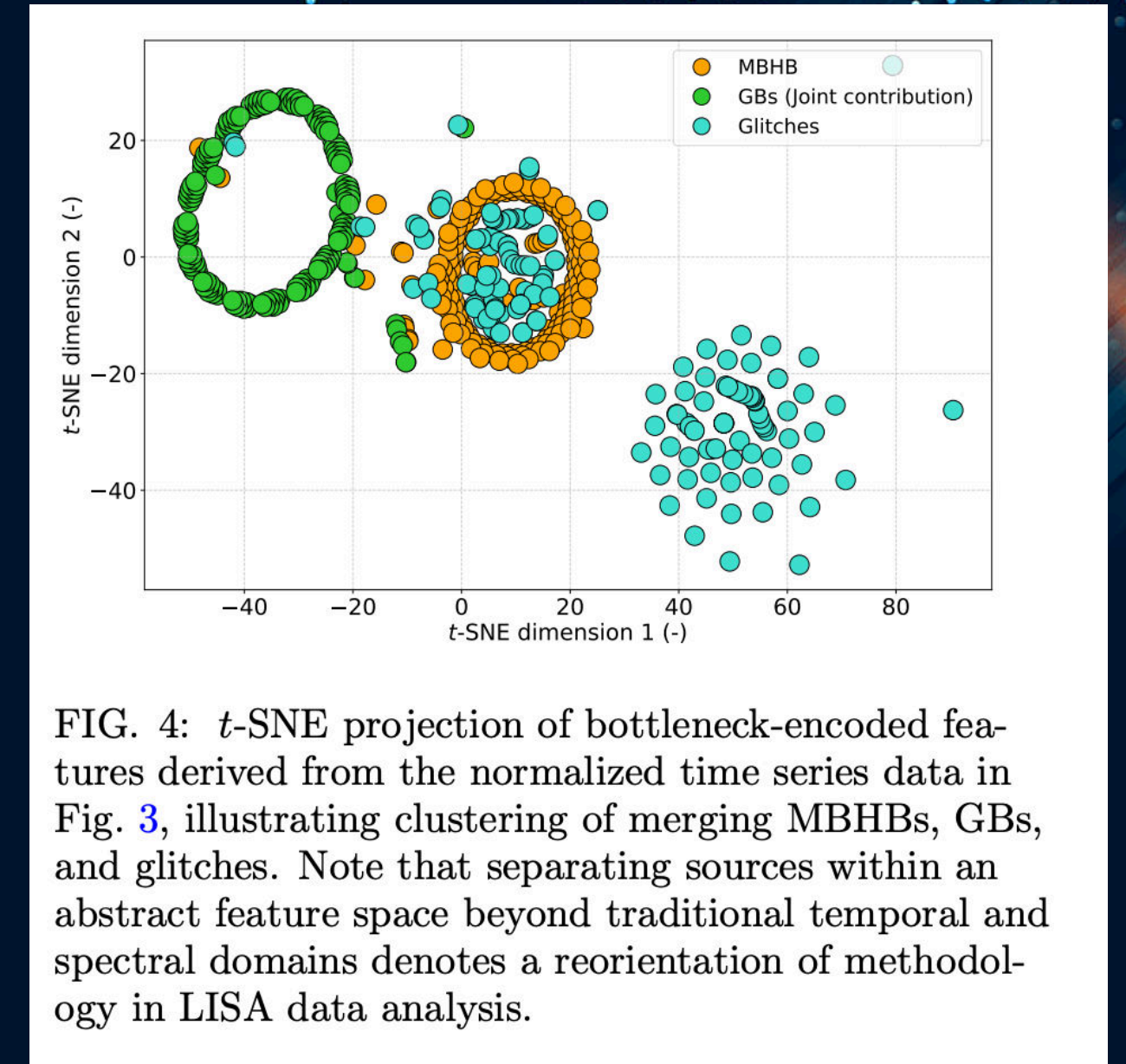
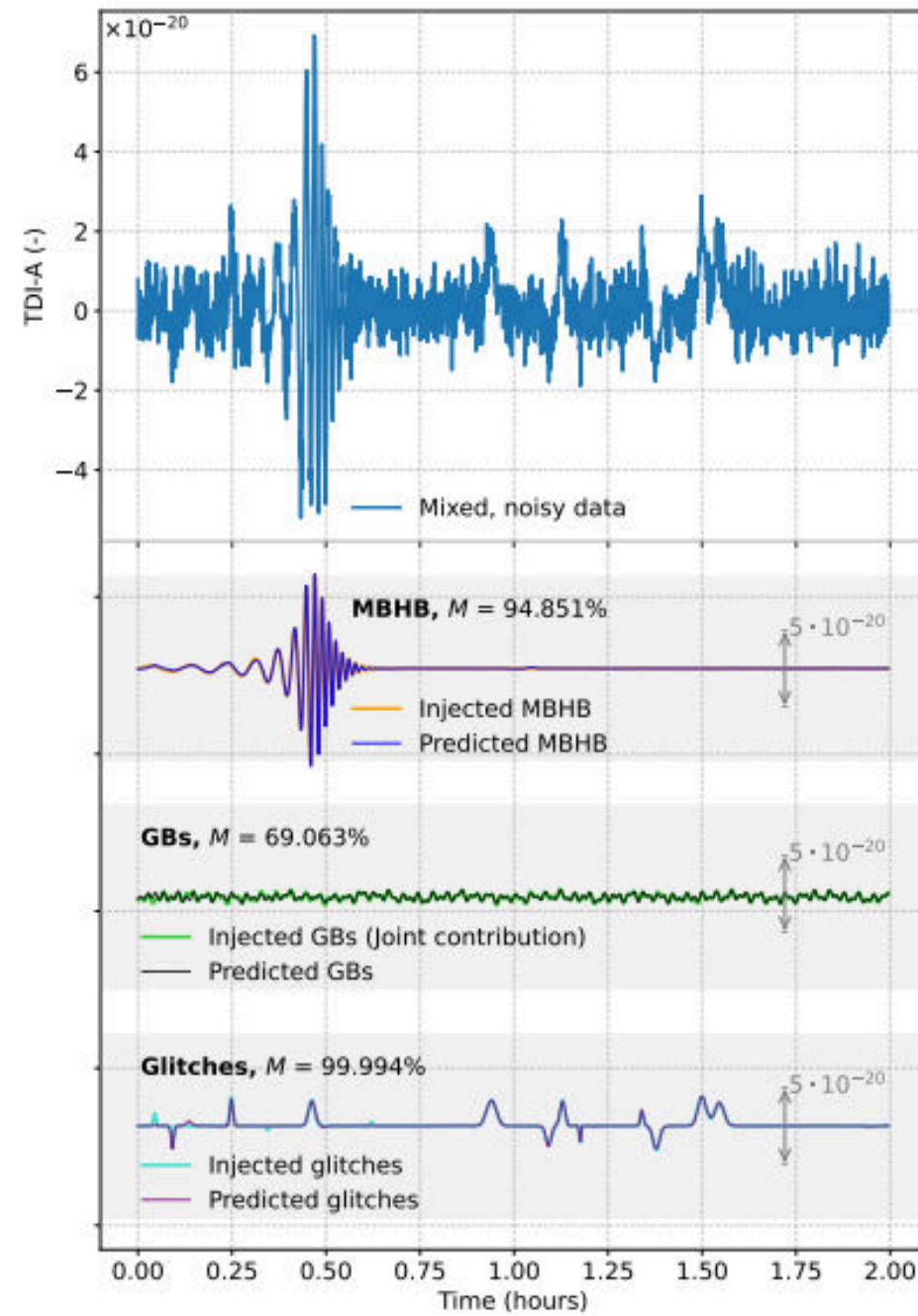


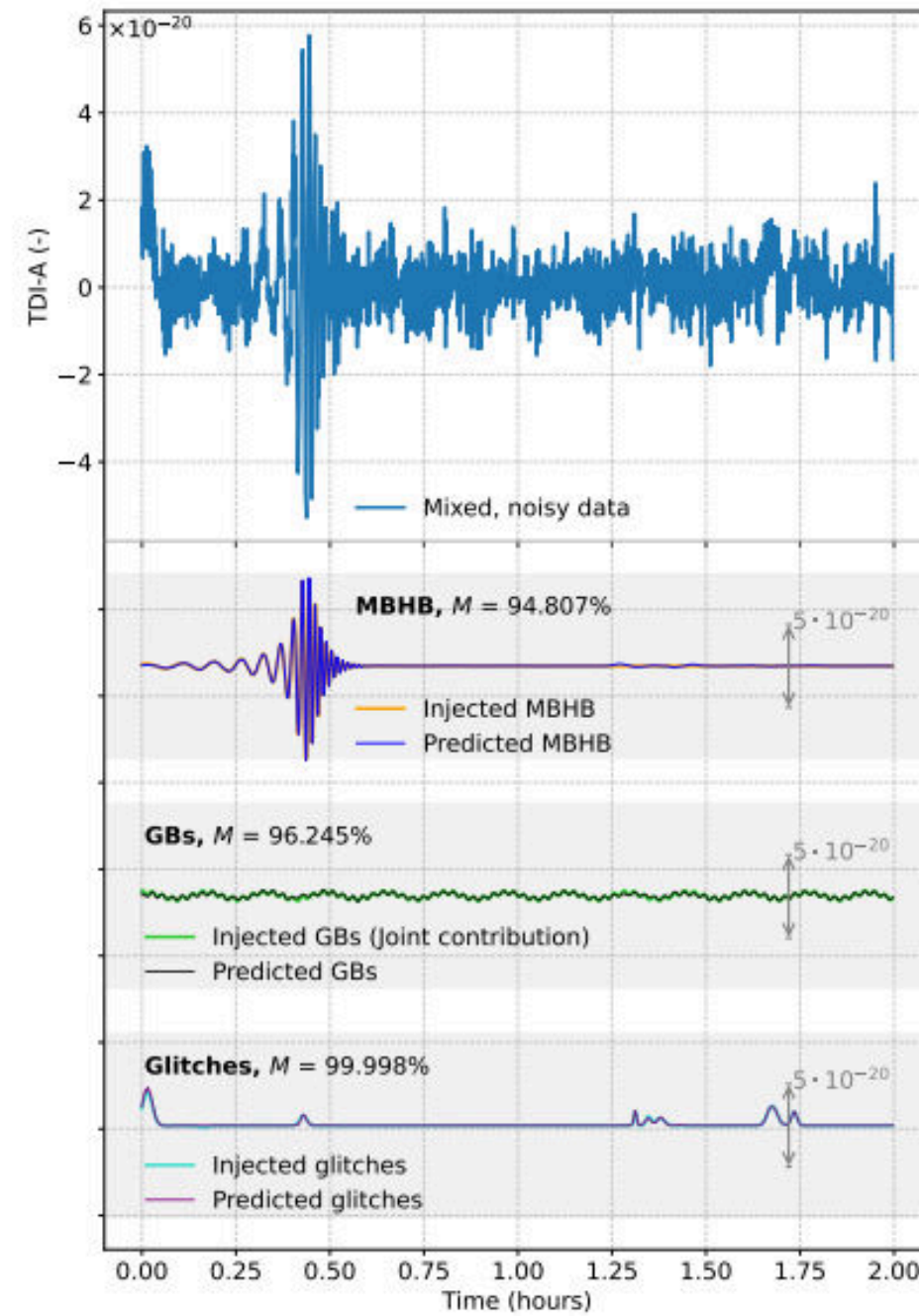
FIG. 4:  $t$ -SNE projection of bottleneck-encoded features derived from the normalized time series data in Fig. 3, illustrating clustering of merging MBHBs, GBs, and glitches. Note that separating sources within an abstract feature space beyond traditional temporal and spectral domains denotes a reorientation of methodology in LISA data analysis.



# Signal denoising & separation



(a)

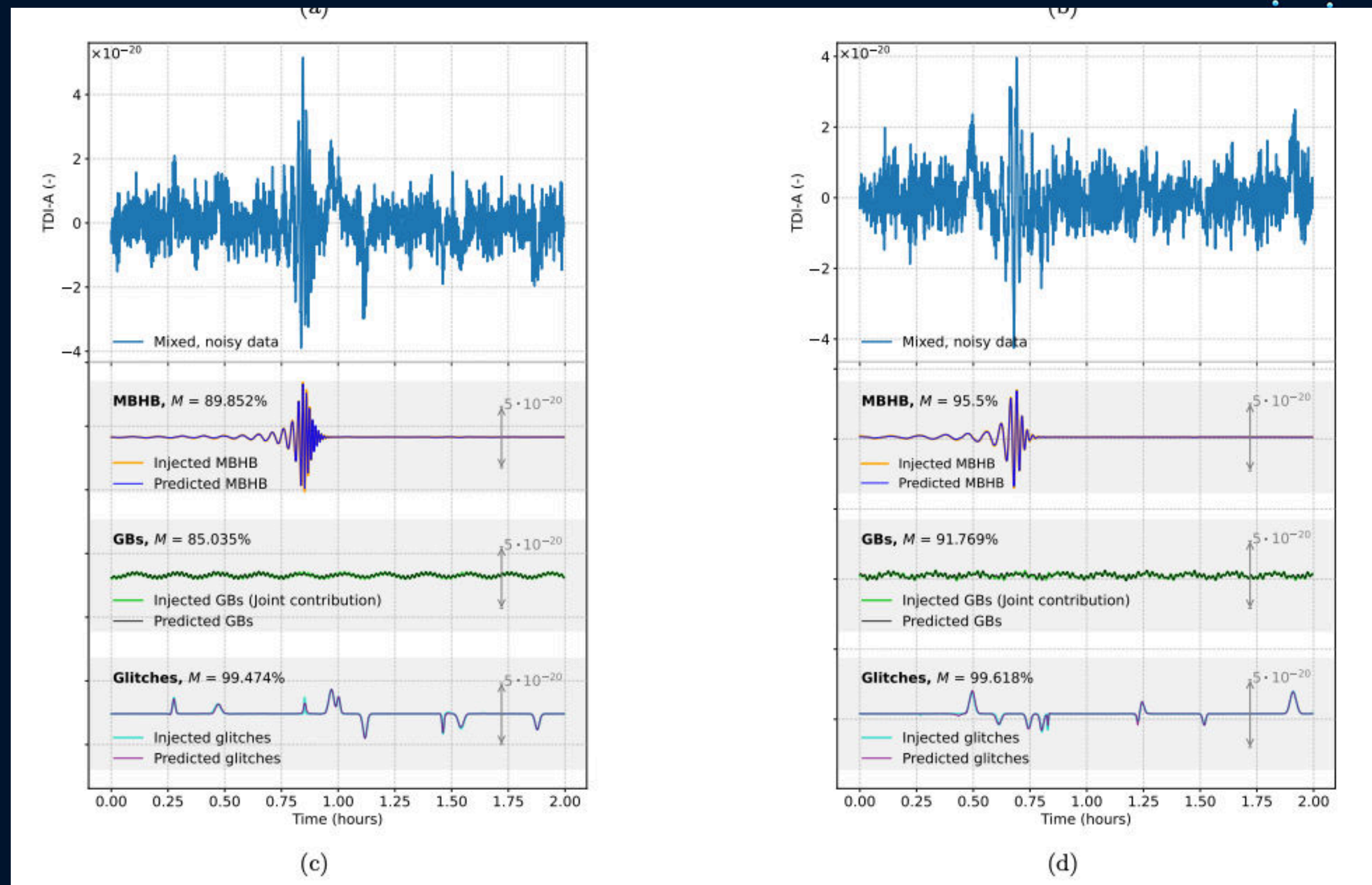


(b)

when glitches overlap with the MBHB merger phase



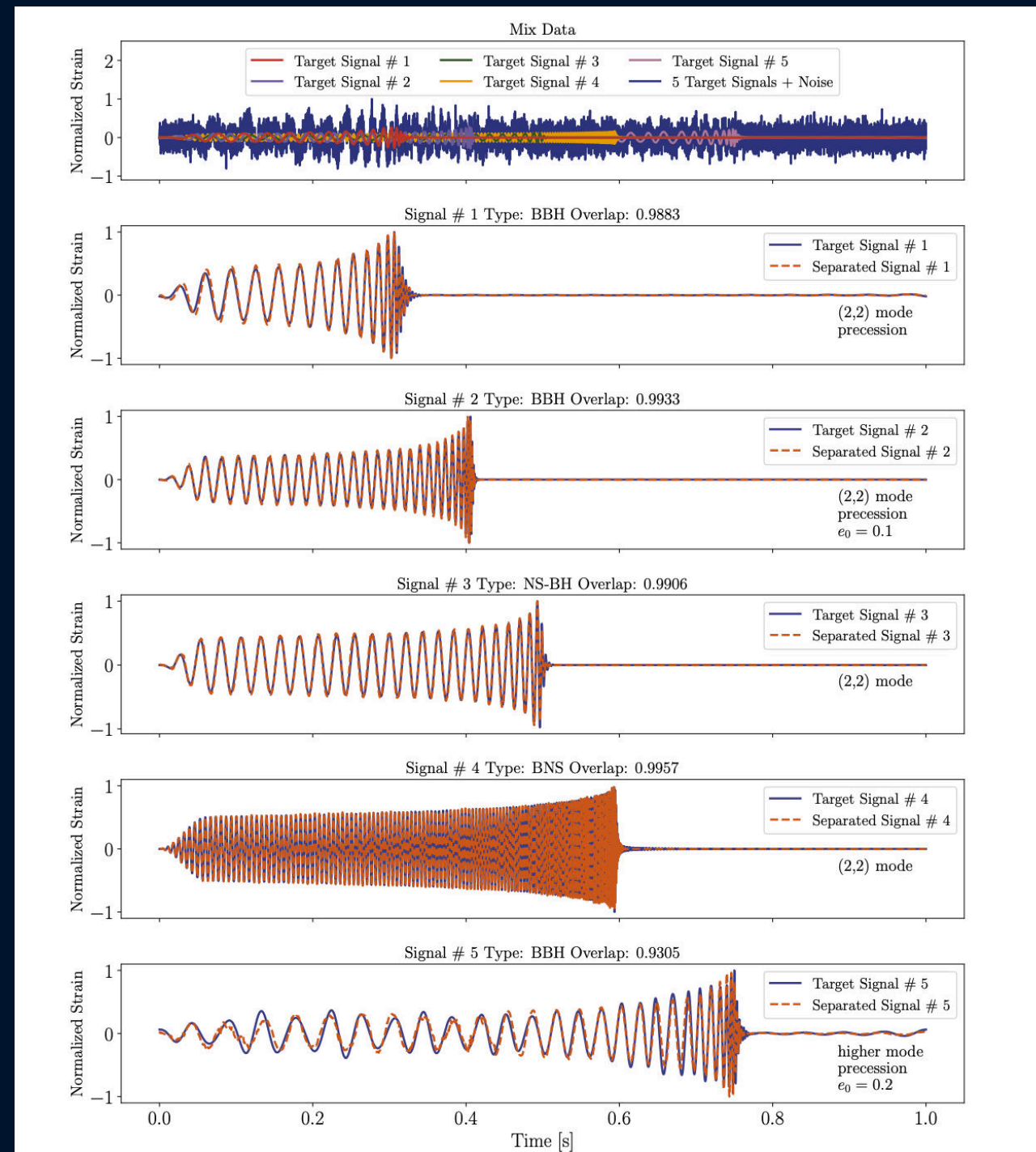
# Signal denoising & separation



when glitches occur during  
the MBHB ringdown

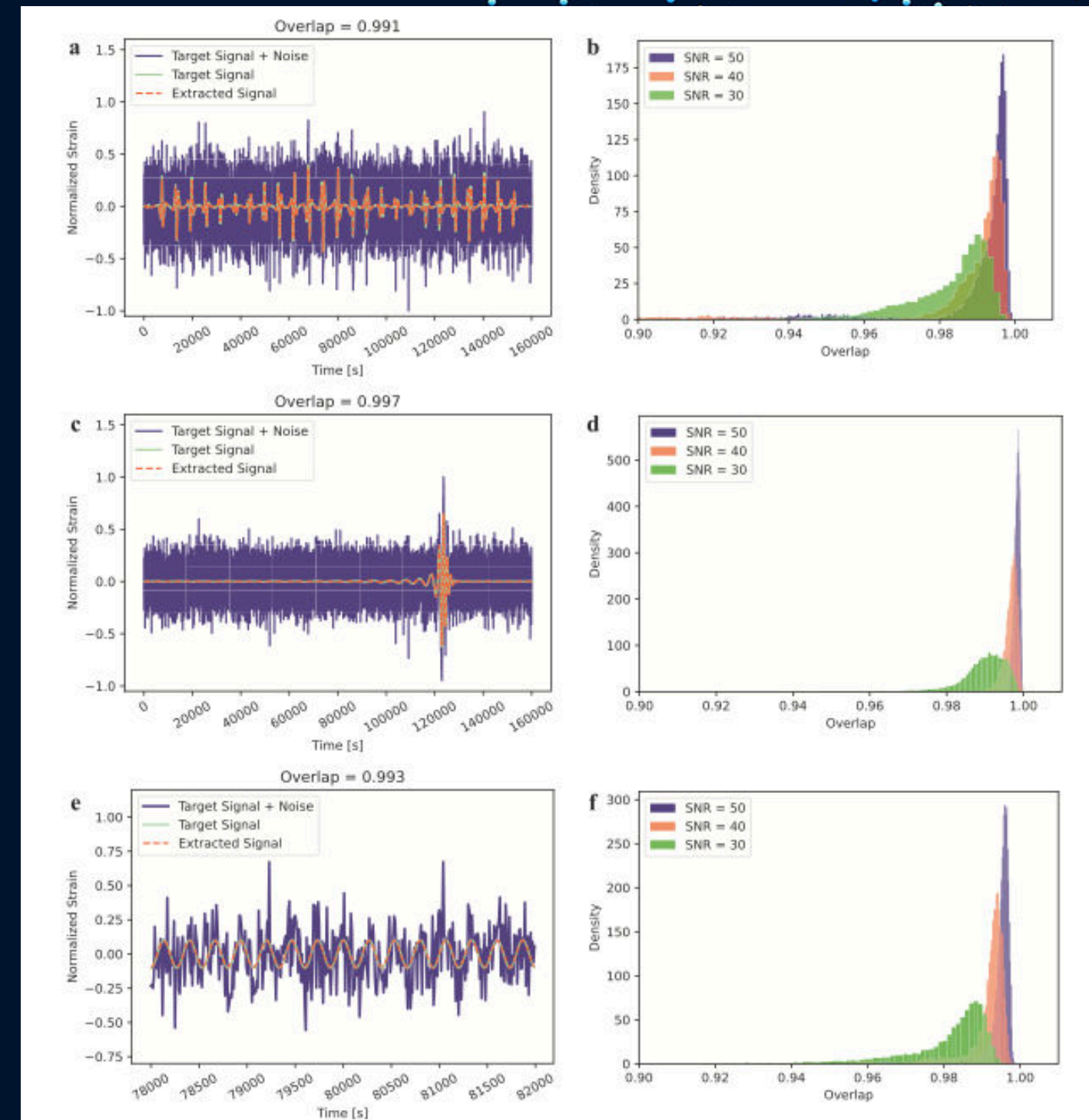


# Signal denoising & separation



**Figure 5. Showcase of signal separation and generalization ability.** The top panel illustrates the mixed data that contains five target signals buried in noise. Other panels display the separated individual signals and target templates, including BBH waveforms with precession, orbital eccentricity, and higher modes, as well as NS-BH and BNS waveforms. The overlaps between the separated and target signals are consistently high, demonstrating the model's effectiveness in separating different types of signals and generalizing to diverse complex waveforms.

<https://arxiv.org/pdf/2412.18259>



**Fig. 4 The signal extraction examples and the overlapping (between the target and extracted signals) distributions for different GW sources. a and b EMRI. c and d MBHB. e and f BWD.** The extracted signal is compared with whitened templates. Only the middle part of the BWD waveform is presented to show the details of the waveform. The overlap between extracted data and waveform templates is shown on the top. The high values indicate the strong performance of our method on signal extraction for different GW sources. Tests on different signal SNRs also show our models' generalization ability.

<https://www.nature.com/articles/s42005-023-01334-6>