
Physics-Informed Discrepancy Decomposition and Robust Astrophysical Inference for GW231123

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Robust astrophysical interpretations from gravitational-wave parameter inference
2 require addressing model-dependent biases. We present a physics-informed frame-
3 work to decompose discrepancies among five waveform models (NRSur7dq4,
4 IMRPhenomXO4a, SEOBNRv5PHM, IMRPhenomXPHM, IMRPhenomTPHM)
5 for GW231123. Our approach combines exploratory metrics (Jensen-Shannon
6 Divergence, Wasserstein distance), high-dimensional analysis with UMAP, and
7 a Physics-Informed Discrepancy Decomposition. This decomposition quantifies
8 divergences in parameter subspaces—mass/distance, effective spin, individual
9 spin/orientation, and remnant properties—linking model differences to physi-
10 cal approximations. We find substantial disagreements in inferred component
11 masses, effective spin, and redshift, with UMAP separating models into distinct
12 clusters. Discrepancy attribution shows individual spin/orientation is most model-
13 dependent due to spin-precession treatments, while remnant properties reflect
14 merger-ringdown modeling. Crucially, no astrophysical parameter for GW231123
15 is robust across all models, as systematic waveform uncertainties exceed statistical
16 errors. Thus, for high-mass, precessing binary black hole mergers, waveform
17 choice dominates inference, limiting firm astrophysical conclusions unless model
18 biases are explicitly accounted for.¹

19

1 Introduction

20 The advent of gravitational-wave (GW) astronomy, beginning with LIGO-Virgo-KAGRA (LVK)
21 detections of binary black hole (BBH) mergers, has transformed our understanding of energetic
22 astrophysical events. GW signals allow us to infer black hole properties, formation channels, and
23 test spacetime under extreme conditions. Yet, such interpretations critically depend on theoretical
24 waveform models, which approximate Einstein’s field equations with varying fidelity—from accurate
25 but costly numerical relativity (NR) simulations to faster semi-analytical and phenomenological
26 models.

27 High-mass, spinning, and precessing BBH systems particularly require approximate models, since
28 NR is too computationally expensive for large-scale inference. Consequently, multiple waveform
29 families exist, differing in efficiency, higher-order mode inclusion, and spin-precession treatment.
30 This diversity introduces model-dependent biases that can dominate statistical errors, especially for
31 challenging events like GW231123. To achieve robust science, it is essential not only to observe

¹This paper, including the idea and the research analysis, was fully generated and written by Denario, a multi-AI agent system. All the input and output files, together with the original paper, can be found in the supplementary material. The Denario code is available in the supplementary material and a YouTube video demonstrating the end-to-end research pipeline with Denario is available in the anonymised YouTube channel at this link.

32 model discrepancies but to identify why they arise and which physical aspects of the source are most
33 sensitive.

34 In this paper, we introduce a physics-informed framework to decompose discrepancies among
35 five models—NRSur7dq4, IMRPhenomXO4a, SEOBNRv5PHM, IMRPhenomXPHM, and IMR-
36 PhenomTPHM [17]. Our key innovation, Physics-Informed Discrepancy Decomposition, defines
37 parameter subspaces (mass/distance, effective spin, individual spin/orientation, remnant properties)
38 and quantifies divergences within them, directly linking differences in inferred parameters to known
39 model approximations such as spin-precession treatment or higher-order mode inclusion.

40 To validate this approach, we first perform exploratory analysis: computing marginal posterior
41 divergences via Jensen-Shannon Divergence (JSD) and 1-Wasserstein distance, and applying Uniform
42 Manifold Approximation and Projection (UMAP) to reveal model clustering in the high-dimensional
43 parameter space. The physics-informed decomposition then quantifies disagreement within each
44 subspace, identifying which GW231123 properties are robust across models and which remain
45 model-sensitive. This provides interpretable insights into astrophysical inference robustness and
46 enables more reliable consensus constraints for GW231123.

47 2 Methods

48 2.1 Data Acquisition and Pre-processing

49 Our analysis initiates with the acquisition and meticulous pre-processing of posterior samples derived
50 from the gravitational-wave event GW231123. These samples, representing the probability distribu-
51 tions of various source parameters, were generated using five distinct gravitational-wave waveform
52 models: NRSur7dq4, IMRPhenomXO4a, SEOBNRv5PHM, IMRPhenomXPHM, and IMRPhenomT-
53 PHM. Each model’s posterior samples were provided as individual CSV files, specifically located at
54 ‘/mnt/ceph/users/anonymous_human/AstroPilot/GW/Iteration1/data/GW231123_NRSur7dq4.csv’,
55 ‘/mnt/ceph/users/anonymous_human/AstroPilot/GW/Iteration1/data/GW231123_IMRPhenomXO4a.csv’,
56 ‘/mnt/ceph/users/anonymous_human/AstroPilot/GW/Iteration1/data/GW231123_SEOBNRv5PHM.csv’,
57 ‘/mnt/ceph/users/anonymous_human/AstroPilot/GW/Iteration1/data/GW231123_IMRPhenomXPHM.csv’,
58 and ‘/mnt/ceph/users/anonymous_human/AstroPilot/GW/Iteration1/data/GW231123_IMRPhenomTPHM.csv’.
59 Upon loading, each CSV file was parsed into a separate pandas DataFrame [26]. To facilitate unified
60 analysis while preserving model attribution, a ‘model’ column was appended to each DataFrame,
61 explicitly identifying the waveform model from which the samples originated. These individual
62 DataFrames were then consolidated into a single, master Python dictionary, with model names serving
63 as keys, providing a structured and accessible representation of the entire dataset.
64 A critical step in pre-processing involved thorough data cleaning and verification [16, 19]. This
65 included confirming the consistency of column names across all files, checking for the presence of
66 any ‘NaN’ or missing values (none were found, as expected for posterior samples of this nature) [16],
67 and verifying the sensible range of $\log_{\text{likelihood}}$ values.

68 2.2 Exploratory Data Analysis and Baseline Comparison

69 Prior to undertaking advanced discrepancy decomposition, we performed an extensive exploratory
70 data analysis to establish a baseline understanding of the agreements and disagreements among the
71 five waveform models. This phase provided initial quantitative insights into the parameter inferences
72 for GW231123 [8, 7].

73 2.2.1 Summary Statistics

74 For each of the five waveform models, we computed key summary statistics for the following
75 astrophysical parameters: $mass_1_source$ (primary component mass), $mass_2_source$ (secondary
76 component mass), chi_eff (effective inspiral spin parameter), chi_p (precessing spin parameter),
77 $redshift$, $final_mass_source$ (remnant black hole mass), and $final_spin$ (remnant black hole
78 spin). Specifically, we calculated the median and the 90% credible interval (defined by the 5th and
79 95th percentiles) for the 1D marginal posterior distribution of each parameter. These statistics were
80 compiled into a comprehensive table, offering an immediate, quantitative overview of the central
81 tendencies and uncertainties predicted by each model.

82 **2.2.2 Pairwise Statistical Divergence**

83 To rigorously quantify the disagreement between the 1D marginal posterior distributions of each
84 parameter across all model pairs, we employed two robust statistical divergence metrics: the Jensen-
85 Shannon Divergence (JSD) [25, 13] and the 1-Wasserstein distance. For each parameter, and for
86 every unique pair of waveform models, the following procedure was applied:

- 87 1. The 1D marginal posterior samples for the given parameter from each model were extracted.
- 88 2. A Kernel Density Estimator (KDE) was used to estimate the Probability Density Function
89 (PDF) for each set of samples. A common, optimized bandwidth (e.g., determined by
90 Scott's rule or Silverman's rule) was applied across all PDFs for a given parameter to ensure
91 consistent smoothing.
- 92 3. The JSD was calculated between the estimated PDFs of the two models. The JSD is
93 a symmetric and finite measure of the similarity between two probability distributions,
94 ranging from 0 (identical distributions) to 1 (maximally divergent distributions), and is based
95 on the Kullback-Leibler divergence.
- 96 4. The 1-Wasserstein distance (also known as Earth Mover's Distance) was computed between
97 the empirical distributions of the two models. This metric quantifies the minimum cost of
98 transforming one distribution into the other, effectively measuring the "distance" between
99 probability distributions.

100 This process yielded two 5x5 symmetric matrices for each key astrophysical parameter, one for JSD
101 values and one for 1-Wasserstein distances. These matrices served as a quantitative baseline for
102 understanding the degree of agreement or disagreement between models on a parameter-by-parameter
103 basis, highlighting where significant univariate discrepancies first emerge.

104 **2.3 High-Dimensional Degeneracy and Discrepancy Analysis**

105 To gain a holistic understanding of how the different waveform models populate the high-dimensional
106 parameter space and to visualize complex degeneracies, we employed Uniform Manifold Approxima-
107 tion and Projection (UMAP) [11, 29].

108 **2.3.1 Data Preparation for UMAP**

109 All posterior samples from the five waveform models were combined into a single, large DataFrame.
110 This consolidated dataset included all 13 physical parameters typically inferred for binary black hole
111 mergers. To ensure that parameters with differing scales did not disproportionately influence the
112 dimensionality reduction, all parameter columns were standardized using z-scoring (subtracting the
113 mean and dividing by the standard deviation across the combined dataset). The 'model' column was
114 retained to allow for post-projection attribution and analysis.

115 **2.3.2 Uniform Manifold Approximation and Projection (UMAP)**

116 The UMAP algorithm was applied to the standardized, high-dimensional parameter space. UMAP is a
117 non-linear dimensionality reduction technique that constructs a high-dimensional graph representing
118 the data's topological structure and then optimizes a low-dimensional graph to be as structurally
119 similar as possible [31]. Our primary goal was to project the high-dimensional parameter space
120 (encompassing all 13 physical parameters) down to a 2D space, thereby enabling intuitive visualization
121 of the complex, non-linear relationships and degeneracies inherent in the posterior distributions
122 [30, 31].

123 We utilized the 'umap-learn' library for this implementation. Key hyperparameters, *n_neighbors*
124 (controlling the balance between local and global structure preservation) and *min_dist* (controlling
125 how tightly points are packed together), were tuned to optimize the embedding quality [18, 18].
126 Initial values of *n_neighbors* = 50 and *min_dist* = 0.1 were used as a starting point [32, 14], with
127 iterative adjustments made to achieve a robust representation that captures both the local clustering
128 and global separation of the data [20]. The output of the UMAP transformation was a set of 2D
129 coordinates (*UMAP_1*, *UMAP_2*) for each posterior sample, representing its position in the
130 learned low-dimensional manifold.

131 **2.3.3 Analysis of UMAP Embedding**

132 The generated 2D UMAP embedding provided a powerful visual and analytical tool to assess the high-
133 dimensional discrepancies [24]. By filtering the UMAP coordinates by their associated 'model' label,
134 we could qualitatively and quantitatively examine how the posterior samples for each waveform model
135 occupy and cluster within this reduced space [10, 4]. We specifically investigated whether the point
136 clouds corresponding to different models exhibited systematic shifts, changes in overall shape, or
137 differences in density concentration. For instance, we analyzed if models with fundamentally different
138 physical approximations, such as IMRPhenomXPHM (which includes a "twisting-up" precession
139 formalism) and NRSur7dq4 (a numerical relativity surrogate), showed distinct, non-overlapping
140 regions in the UMAP space [12], indicating significant high-dimensional disagreements.

141 **2.4 Physics-Informed Discrepancy Decomposition**

142 The core of our methodology lies in the Physics-Informed Discrepancy Decomposition [9], which
143 systematically dissects the overall model disagreements and attributes them to specific physical
144 effects [9] and the corresponding approximations within the waveform models [22]. This approach
145 goes beyond global comparisons by focusing on physically motivated parameter subspaces [9].

146 **2.4.1 Definition of Physical Parameter Subspaces**

147 Based on our understanding of binary black hole physics and the known characteristics and approxi-
148 mation schemes of the waveform models [21, 23, 15], we meticulously defined four distinct parameter
149 subspaces. These subspaces were designed to isolate specific physical aspects of the binary merger
150 that are known to be treated differently across waveform models [15]. For each subspace, we created
151 subsets of the posterior data, containing only the relevant parameters.

- 152 **1. Mass & Distance Subspace:** This subspace includes (*mass_1_source*, *mass_2_source*,
153 *redshift*). These parameters are fundamental to the overall amplitude and frequency
154 evolution of the gravitational-wave signal. Discrepancies in this subspace can often be
155 attributed to differences in the leading-order inspiral dynamics or the calibration against
156 astrophysical priors.
- 157 **2. Effective Spin Subspace:** Comprising (*chi_eff*, *chi_p*), this subspace captures the domi-
158 nant, orbit-averaged effects of spin. *chi_eff* primarily influences the inspiral rate, while
159 *chi_p* quantifies the strength of orbital plane precession. Disagreements here reflect how
160 models approximate the average spin effects throughout the inspiral.
- 161 **3. Individual Spin & Orientation Subspace:** This is a high-dimensional subspace consisting
162 of (*a_1*, *a_2*, *cos_tilt_1*, *cos_tilt_2*, *cos_theta_jn*, *phi_jl*). These parameters describe
163 the detailed magnitudes and orientations of the individual black hole spins, as well as the ori-
164 entation of the binary's orbital angular momentum relative to the line of sight. This subspace
165 is particularly sensitive to the treatment of spin precession, including the full precessional
166 dynamics (as in NRSur7dq4 and SEOBNRv5PHM) versus simplified "twisting-up" for-
167 malisms (as in IMRPhenomXPHM and IMRPhenomTPHM). Significant discrepancies in
168 this subspace directly indicate differences in how models handle the complex interplay of
169 spins and orbital dynamics.
- 170 **4. Remnant Properties Subspace:** This subspace includes (*final_mass_source*,
171 *final_spin*). These parameters represent the predicted properties of the final black hole
172 formed after the merger. They are highly sensitive to the modeling of the merger-ringdown
173 phase of the waveform, as well as the accurate inclusion of higher-order waveform modes,
174 which become more prominent during this phase.

175 **2.4.2 Quantifying Subspace-Specific Discrepancies**

176 For each of the four defined physical subspaces, and for every pairwise combination of the five
177 waveform models, we quantified the multi-dimensional disagreement using the multi-dimensional
178 Jensen-Shannon Divergence (JSD) [25, 13]. The procedure for computing multi-dimensional JSD for
179 a given subspace between two models (e.g., Model A and Model B) was as follows:

- 180 1. The posterior samples for the parameters within the specific subspace were extracted for
181 both Model A and Model B.

- 182 2. A multi-dimensional Kernel Density Estimator (KDE) was employed to estimate the joint
 183 PDF for each model's samples within that subspace. This involves estimating the probability
 184 density across the entire multi-dimensional space spanned by the subspace parameters.
 185 3. The multi-dimensional JSD was then computed between the two estimated joint PDFs. This
 186 metric provides a single scalar value quantifying the overall divergence of the two models'
 187 posterior distributions within that specific physical subspace.

188 This process resulted in four separate 5x5 discrepancy matrices [28], one for each physical subspace.
 189 Each matrix element represented the multi-dimensional JSD [1, 6] between a pair of models within
 190 that particular subspace, thereby providing a targeted measure of disagreement.

191 **2.4.3 Correlation of Discrepancies with Model Physics**

192 The resulting discrepancy matrices from the physics-informed decomposition were critically analyzed
 193 to establish direct links between the magnitude of the observed discrepancies and the known physical
 194 differences in the underlying waveform models.

195 For instance, we specifically compared the JSD values in the 'Individual Spin & Orientation' matrix
 196 with those in the 'Mass & Distance' matrix. We hypothesized that models with fundamentally
 197 different treatments of spin precession (e.g., IMRPhenomXPHM versus NRSur7dq4) would exhibit
 198 significantly larger JSD values in the highly sensitive spin and orientation subspace compared to the
 199 more universally agreed-upon mass and distance subspace. Similarly, we examined the 'Remnant
 200 Properties' discrepancy matrix, anticipating that models incorporating a more complete treatment
 201 of higher-order modes (such as SEOBNRv5PHM and IMRPhenomXPHM) would show greater
 202 consistency among themselves, while displaying larger divergences with models that have a less
 203 comprehensive representation of the merger-ringdown phase, like IMRPhenomXO4a. This systematic
 204 correlation allowed us to attribute discrepancies to specific physical approximations within the models,
 205 moving beyond mere observation of disagreement to understanding its underlying causes.

206 **2.5 Robust Astrophysical Inference**

207 The final stage of our analysis involved synthesizing the findings from the exploratory data analysis,
 208 high-dimensional embedding, and physics-informed decomposition to derive robust astrophysical
 209 constraints for GW231123 [33, 5, 27].

210 **2.5.1 Identification of Robustly Constrained Parameters**

211 A key objective was to identify which astrophysical parameters for GW231123 are robustly con-
 212 strained across all five waveform models, meaning their inferred posterior distributions show high
 213 consistency regardless of the model choice [2, 3]. A parameter was deemed "robust" if the maximum
 214 pairwise Jensen-Shannon Divergence (JSD) and 1-Wasserstein distance values among all model pairs
 215 (as calculated in Section 2.2) fell below a pre-defined threshold (e.g., $JSD < 0.01$). Furthermore,
 216 strong overlap in the medians and 90% credible intervals across all models, as observed in the
 217 summary statistics, served as an additional indicator of robustness [2].

218 **2.5.2 Identification of Model-Dependent Parameters**

219 Conversely, parameters that failed to meet the robustness criteria were classified as "model-
 220 dependent." For these parameters, the systematic uncertainties introduced by waveform model
 221 choice were found to be significant. Crucially, our Physics-Informed Discrepancy Decomposition
 222 (Section 4.3) allowed us to pinpoint the primary physical origin of these discrepancies. For example,
 223 if $\chi_{i,p}$ was identified as model-dependent, the analysis would then attribute this discrepancy to
 224 differing treatments of spin precession between phenomenological and NR-calibrated models, based
 225 on the high JSD values observed in the 'Individual Spin & Orientation' subspace.

226 **2.5.3 Derivation of Consensus Astrophysical Constraints**

227 For those parameters identified as robustly constrained, we derived a final consensus measurement
 228 for GW231123 [7]. This was achieved by combining the posterior samples for that specific parameter
 229 from all five waveform models into a single, aggregated dataset. From this combined distribution, the

230 final consensus median and 90% credible interval were computed, representing our most reliable,
231 model-agnostic measurement for that property of the binary black hole system [7].

232 **2.5.4 Final Results Compilation**

233 The comprehensive findings were compiled into a final summary table. This table explicitly listed all
234 key astrophysical parameters of GW231123. For each parameter, it provided the derived consensus
235 median and 90% credible interval if the parameter was deemed robust.

236 If a parameter was classified as model-dependent, the table reported the range of medians observed
237 across the different models instead of a single consensus value, clearly marking it as such. An
238 additional column provided a concise statement on whether the parameter constraint was 'Robust' or
239 'Model-Dependent', along with a brief, physics-informed note explaining the origin of any significant
240 model dependency, directly linking back to the insights gained from the discrepancy decomposition.
241 This structured presentation allowed for a clear and interpretable assessment of the astrophysical
242 inferences for GW231123.

243 **3 Results**

244 **3.1 Baseline comparison: Significant divergence in key physical parameters**

245 Our initial exploratory data analysis, utilizing summary statistics and pairwise statistical divergence
246 metrics as outlined in Section 2.2, immediately revealed substantial disagreements among the five
247 waveform models regarding the inferred astrophysical parameters for GW231123. As summarized in
248 Table 1 and visually presented through one-dimensional marginal posterior distributions in Figure 1,
249 key source parameters exhibit significant model-dependent variations.

250 The most pronounced discrepancy, evident in both Table 1 and Figure 1, is observed in the component
251 masses, particularly for `mass_2_source`. While NRSur7dq4, SEOBNRv5PHM, and IMRPhenomTPHM infer a relatively symmetric binary system with `mass_2_source` medians ranging from
252 $110.04 M_{\odot}$ to $111.10 M_{\odot}$, IMRPhenomXO4a predicts a significantly more asymmetric configuration,
253 with a median `mass_2_source` of only $55.08 M_{\odot}$. IMRPhenomXPHM also infers a lower secondary
254 mass ($93.33 M_{\odot}$) compared to the first group, further highlighting model-dependent variations. This
255 fundamental disagreement in the mass ratio propagates to other inferred parameters, such as the
256 effective inspiral spin parameter (`chi_eff`) and `redshift`.

257 For `chi_eff`, the inferred median values span a considerable range, from a near-zero value of 0.04 for
258 IMRPhenomXPHM to a significantly positive 0.44 for SEOBNRv5PHM and IMRPhenomTPHM, as
259 shown in Table 1 and visually confirmed by the distinct posterior peaks in Figure 1. Such a wide range
260 has profound implications for understanding the astrophysical formation channels of GW231123, as
261 `chi_eff` is a key indicator of the binary's spin alignment with the orbital angular momentum. In
262 contrast, the precessing spin parameter (`chi_p`) shows a comparatively smaller spread in median
263 values (from 0.73 to 0.82), suggesting that while the magnitude of precession is consistently inferred
264 to be high, its detailed influence on other parameters varies.

265 These disagreements are quantitatively supported by the pairwise Jensen-Shannon Divergence (JSD)
266 and 1-Wasserstein distance metrics, calculated as described in Section 2.2. For instance, JSD values
267 between certain model pairs for `mass_2_source` and `redshift` frequently exceed 0.6, indicating
268 near-complete non-overlap of the 1D marginal posterior distributions, as is clearly visible in Figure
269 1. For `redshift`, IMRPhenomXPHM consistently places the source at a much closer distance
270 (median 0.17), while IMRPhenomXO4a infers a significantly more distant source (median 0.58),
271 with other models falling in between. This initial assessment underscores that the choice of waveform
272 model introduces substantial systematic uncertainties that cannot be overlooked in astrophysical
273 interpretations.

275 **3.2 High-dimensional degeneracy and model clustering**

276 To gain a more comprehensive understanding of how the waveform models populate the full,
277 high-dimensional parameter space, we employed Uniform Manifold Approximation and Projection
278 (UMAP), as detailed in Section 2.3. The 2D UMAP embeddings, generated from the 13-dimensional

279 parameter space and shown in Figure 2 and Figure 3, provide a powerful visualization of the complex
280 degeneracies and discrepancies.

281 The UMAP projection, as depicted in Figure 2 and Figure 3, clearly reveals a structured separation
282 of the models into distinct clusters. This indicates that the discrepancies are not merely isolated to
283 individual parameters but are inherent to the correlated, high-dimensional posterior distributions. The
284 models coalesce into three primary groups:

- 285 1. **A Core Cluster:** Comprising NRSur7dq4, SEOBNRv5PHM, and IMRPhenomTPHM.
286 These models occupy a contiguous region in the UMAP embedding, suggesting a higher
287 degree of consistency in their high-dimensional parameter inferences.
- 288 2. **An Isolated Cluster (IMRPhenomXO4a):** This model forms a distinct, separate cluster,
289 indicating significant divergence from all other models in the overall parameter space.
- 290 3. **A Second Isolated Cluster (IMRPhenomXPHM):** This model also forms a unique cluster,
291 located in a region of the UMAP space far from the other models.

292 Table 2 provides the UMAP centroid coordinates for each model, quantitatively illustrating their
293 separation in the learned low-dimensional manifold. IMRPhenomXPHM is positioned at $UMAP_1 \approx -3.86$, while IMRPhenomXO4a is at $UMAP_1 \approx 11.42$, confirming their extreme separation from
294 the core cluster which is centered around $UMAP_1$ values closer to 0 – 3.

296 This clustering is physically meaningful. The two most separated models, IMRPhenomXO4a and
297 IMRPhenomXPHM, are both frequency-domain phenomenological models, but they incorporate
298 different physical approximations, particularly in their treatment of higher-order modes and spin
299 precession. For instance, IMRPhenomXPHM employs a "twisting-up" formalism for precession,
300 which differs from the more complete dynamical evolution captured by numerical relativity (NR)
301 surrogates like NRSur7dq4 and effective-one-body (EOB) models like SEOBNRv5PHM. The relative
302 agreement within the core cluster suggests that for a high-mass, potentially precessing system like
303 GW231123, the NR-calibrated and EOB-based time-domain models, along with the time-domain
304 phenomenological model IMRPhenomTPHM, provide more consistent descriptions of the underlying
305 physical dynamics. The UMAP analysis thus serves as a powerful diagnostic tool, demonstrating that
306 waveform model choice fundamentally alters the inferred parameter space for GW231123.

307 3.3 Physics-informed discrepancy decomposition

308 To link high-dimensional disagreements to physical effects and approximations, we performed a
309 physics-informed discrepancy decomposition (Section 2.4). Multi-dimensional Jensen-Shannon
310 Divergence (JSD) was quantified between model pairs in four parameter subspaces: Mass & Distance,
311 Effective Spin, Individual Spin & Orientation, and Remnant Properties, with results shown in Figure
312 4.

313 3.3.1 Mass & distance subspace

314 Parameters `mass_1_source`, `mass_2_source`, and `redshift` show very high JSD values (often
315 > 0.6), confirming strong model disagreement on intrinsic masses and source distance (Figure
316 4, top-left). Degeneracies with spin and orientation drive large shifts, showing even basic source
317 properties are not robust without modeling systematics.

318 3.3.2 Effective spin subspace

319 Effective spin parameters (`chi_eff`, `chi_p`) also show large discrepancies (Figure 4, top-right).
320 IMRPhenomXPHM vs. IMRPhenomTPHM diverge strongly ($JSD = 0.636$), reflecting different
321 spin-orbit treatments. By contrast, SEOBNRv5PHM and IMRPhenomTPHM agree well ($JSD = 0.043$), indicating consistent orbit-averaged spin modeling despite distinct paradigms.

323 3.3.3 Individual spin & orientation subspace

324 The 6-D space (`a1`, `a2`, `cos_tilt_1`, `cos_tilt_2`, `cos_theta_jn`, `phi_j1`) shows the strongest
325 disagreements, with many JSD values near the 0.693 maximum (Figure 4, bottom-left). This reflects
326 model-dependent spin-precession treatments: simplified "twisting-up" models (IMRPhenomXPHM,

327 IMRPhenomXO4a) yield different posteriors than fully dynamical precession models (NRSur7dq4,
328 SEOBNRv5PHM), making spin configuration highly uncertain.

329 **3.3.4 Remnant properties subspace**

330 Final black hole properties (`final_mass_source`, `final_spin`) are also model-dependent (Figure
331 4, bottom-right). IMRPhenomXPHM predicts lower final spin (median 0.71) than others (0.81–0.89,
332 Table 1), reflecting differences in merger-ringdown modeling and numerical relativity calibration.
333 Inclusion of higher-order modes is critical. SEOBNRv5PHM and IMRPhenomTPHM again agree
334 well ($JSD = 0.051$).

335 **3.4 Robust astrophysical inference for GW231123**

336 Combining all analyses, robustness was defined (Section 2.5) as pairwise $JSD < 0.05$ and median
337 spread $< 10\%$. Table 3 shows that *no parameter for GW231123 is robust*; model differences prevent
338 consensus.

339 For high-mass, precessing BBHs like GW231123, short merger–ringdown–dominated signals mean
340 waveform-choice systematics match or exceed statistical errors. Wide spreads—`mass_2_source`
341 55.1–111.1, M_\odot , `chi_eff` 0.04–0.44—make distinguishing formation channels impossible without
342 systematic treatment. For GW231123, waveform choice dominates interpretation, precluding firm
343 astrophysical conclusions.

344 **4 Conclusions**

345 **4.1 Problem Statement and Our Approach**

346 Astrophysical interpretation of GW events, especially complex BBH mergers like GW231123, is
347 hindered by model-dependent biases from approximate waveform models. These models trade
348 physical fidelity for efficiency, introducing systematic uncertainties larger than statistical errors.
349 This paper addressed the issue with a physics-informed framework that decomposes and attributes
350 discrepancies among waveform models, quantifying multi-dimensional divergences in parameter
351 subspaces to link model differences to specific physical approximations.

352 **4.2 Summary of Findings**

353 Our comprehensive analysis of GW231123, utilizing five distinct waveform models (NRSur7dq4,
354 IMRPhenomXO4a, SEOBNRv5PHM, IMRPhenomXPHM, IMRPhenomTPHM), yielded several
355 key findings:

- 356 **1. Significant Baseline Disagreements:** Initial exploratory data analysis revealed substantial
357 discrepancies in 1D marginal posterior distributions for key astrophysical parameters, most
358 notably for component masses (especially `mass_2_source`), effective inspiral spin (`chi_eff`),
359 and redshift. The Jensen-Shannon Divergence (JSD) and 1-Wasserstein distance metrics
360 frequently indicated near-complete non-overlap between certain model pairs.
- 361 **2. High-Dimensional Model Clustering:** Uniform Manifold Approximation and Projec-
362 tion (UMAP) confirmed that these discrepancies are not isolated but permeate the high-
363 dimensional parameter space. The UMAP embedding clearly separated the models into
364 distinct clusters, with NRSur7dq4, SEOBNRv5PHM, and IMRPhenomTPHM forming
365 a core cluster, while IMRPhenomXO4a and IMRPhenomXPHM occupied significantly
366 isolated regions. This clustering directly reflects fundamental differences in how these
367 models describe the underlying physical dynamics of GW231123.
- 368 **3. Physics-Informed Discrepancy Attribution:** Our core Physics-Informed Discrepancy
369 Decomposition successfully attributed these model differences to specific physical approxi-
370 mations:
 - 371 • The *Mass & Distance subspace* showed high JSD values, indicating that even funda-
372 mental source properties like masses and redshift are strongly degenerate with and
373 sensitive to the overall waveform modeling.

- 374 • The *Effective Spin subspace* exhibited substantial disagreements, particularly between
 375 IMRPhenomXPHM and IMRPhenomTPHM, highlighting differing treatments of spin-
 376 orbit coupling.
 377 • The *Individual Spin & Orientation subspace* revealed the most severe model depen-
 378 dence, with JSD values approaching maximum divergence. This is a direct consequence
 379 of the varying formalisms for spin precession (e.g., full dynamical precession versus
 380 simplified "twisting-up" approximations) employed by the models.
 381 • The *Remnant Properties subspace* also showed significant model dependence, sensitive
 382 to the modeling of the merger-ringdown phase and the inclusion of higher-order
 383 waveform modes, which are crucial for accurately predicting the final black hole's
 384 mass and spin.

385 **4. Lack of Robust Constraints:** Crucially, our analysis concluded that *no key astrophysical*
 386 *parameter for GW231123 is robustly constrained across all five waveform models.* The sys-
 387 *tematic uncertainties introduced by waveform model choice consistently exceeded statistical*
 388 *uncertainties for this event.*

389 4.3 Implications for Astrophysical Inference

390 This work shows that for high-mass, potentially precessing binary black hole mergers like GW231123,
 391 the waveform model choice is not minor but central to interpretation. Inferred values of component
 392 masses, effective spin, and redshift vary widely, affecting conclusions about the source's nature and
 393 history. For example, the mass ratio spread (mass_2_source from 55.1, M_{\odot} to 111.1, M_{\odot}) can lead
 394 to very different origin scenarios.

395 Our decomposition clarifies that spin precession treatment and merger-ringdown modeling drive these
 396 model-dependent biases. This highlights the need for waveform models that capture spin precession
 397 and higher-order modes. Reliable astrophysical inference will require either consistent models
 398 across key parameter spaces or systematic uncertainty methods that account for model discrepancies.
 399 Without this, conclusions about extreme GW events remain uncertain.

400 References

- 401 [1] Celal Alagoz. Exploring hierarchical classification performance for time series data: Dissimilarity
 402 measures and classifier comparisons, 2024.
- 403 [2] P. S. Aswathi, William E. East, Nils Siemonsen, Ling Sun, and Dana Jones. Ultralight boson
 404 constraints from gravitational wave observations of spinning binary black holes, 2025.
- 405 [3] Imre Bartos and Zoltan Haiman. Accretion is all you need: Black hole spin alignment in merger
 406 gw231123 indicates accretion pathway, 2025.
- 407 [4] F. Bufano, C. Bordiu, T. Ceccarello, M. Munari, A. Hopkins, A. Ingallinera, P. Leto, S. Loru,
 408 S. Riggi, E. Sciacca, G. Vizzari, A. De Marco, C. S. Buemi, F. Cavallaro, C. Trigilio, and
 409 G. Umana. Sifting the debris: Patterns in the SNR population with unsupervised ML methods,
 410 2024.
- 411 [5] Andrea Caputo, Gabriele Franciolini, and Samuel J. Witte. Superradiance constraints from
 412 gw231123, 2025.
- 413 [6] Prateek Chanda, Saral Sureka, Parth Pratim Chatterjee, Krishnateja Killamsetty, Nikhil Shiv-
 414 akumar Nayak, and Ganesh Ramakrishnan. Learning what matters: Probabilistic task selection
 415 via mutual information for model finetuning, 2025.
- 416 [7] The LIGO Scientific Collaboration, the Virgo Collaboration, and the KAGRA Collaboration.
 417 Gw231123: a binary black hole merger with total mass 190-265 m_{\odot} , 2025.
- 418 [8] Iuliu Cuceu, Marie Anne Bizouard, Nelson Christensen, and Mairi Sakellariadou. Gw231123:
 419 Binary black hole merger or cosmic string?, 2025.
- 420 [9] Josef Dick, Seungchan Ko, Kassem Mustapha, and Sanghyeon Park. Locking-free training of
 421 physics-informed neural network for solving nearly incompressible elasticity equations, 2025.

- 422 [10] Dimple, K. Misra, and K. G. Arun. Evidence for two distinct populations of kilonova-associated
423 gamma ray bursts, 2023.
- 424 [11] Benyamin Ghojogh, Ali Ghodsi, Fakhri Karray, and Mark Crowley. Uniform manifold approxi-
425 mation and projection (UMAP) and its variants: Tutorial and survey, 2021.
- 426 [12] Niklas Houba. Deep source separation of overlapping gravitational-wave signals and non-
427 stationary noise artifacts, 2025.
- 428 [13] Jhoan K. Hoyos-Osorio and Luis G. Sanchez-Giraldo. The representation jensen-shannon
429 divergence, 2024.
- 430 [14] Myeongwon Jung, Takanori Fujiwara, and Jaemin Jo. Ghostumap2: Measuring and analyzing
431 (r,d)-stability of UMAP, 2025.
- 432 [15] Veome Kapil, Luca Reali, Roberto Cotesta, and Emanuele Berti. Systematic bias from waveform
433 modeling for binary black hole populations in next-generation gravitational wave detectors,
434 2024.
- 435 [16] Samiya Khan, Xiufeng Liu, and Mansaf Alam. A spark ML driven preprocessing approach for
436 deep learning based scholarly data applications, 2019.
- 437 [17] Yin-Jie Li, Shan-Peng Tang, Ling-Qin Xue, and Yi-Zhong Fan. Gw231123: a product of
438 successive mergers from ~ 10 stellar-mass black holes, 2025.
- 439 [18] Yin-Ting Liao, Hengrui Luo, and Anna Ma. Efficient and robust bayesian selection of hyperpa-
440 rameters in dimension reduction for visualization, 2023.
- 441 [19] Mateus Guimarães Lima, Antony Carvalho, João Gabriel Álvares, Clayton Escouper das
442 Chagas, and Ronaldo Ribeiro Goldschmidt. Impacts of data preprocessing and hyperparameter
443 optimization on the performance of machine learning models applied to intrusion detection
444 systems, 2024.
- 445 [20] Justin Lin and Julia Fukuyama. Calibrating dimension reduction hyperparameters in the
446 presence of noise, 2024.
- 447 [21] Xiaolin Liu, Zhoujian Cao, and Lijing Shao. Upgraded waveform model of eccentric binary
448 black hole based on effective-one-body-numerical-relativity for spin-aligned binary black holes,
449 2023.
- 450 [22] Xinru Mu, Shijun Cheng, and Tariq Alkhalifah. Separationpinn: Physics-informed neural
451 networks for seismic p- and s-wave mode separation, 2025.
- 452 [23] Samanwaya Mukherjee, Khun Sang Phukon, Sayak Datta, and Sukanta Bose. Phenomenological
453 gravitational waveform model of binary black holes incorporating horizon fluxes, 2024.
- 454 [24] Michela Negro, Nicoló Cibrario, Eric Burns, Joshua Wood, Adam Goldstein, and Tito Dal
455 Canton. Prompt GRB recognition through waterfalls and deep learning, 2025.
- 456 [25] Frank Nielsen. On a generalization of the jensen-shannon divergence and the JS-symmetrization
457 of distances relying on abstract means, 2022.
- 458 [26] Shriram Shanbhag and Sridhar Chimalakonda. On the energy consumption of different
459 dataframe processing libraries – an exploratory study, 2022.
- 460 [27] Ataru Tanikawa, Shuai Liu, WeiWei Wu, Michiko S. Fujii, and Long Wang. Gw231123
461 formation from population III stars: Isolated binary evolution, 2025.
- 462 [28] Michail Tsagris. Constrained least squares simplicial-simplicial regression, 2024.
- 463 [29] Brent Vela, Trevor Hastings, and Raymundo Arróyave. Visualizing high entropy alloy spaces:
464 Methods and best practices, 2024.
- 465 [30] Yingfan Wang, Haiyang Huang, Cynthia Rudin, and Yaron Shaposhnik. Understanding how
466 dimension reduction tools work: An empirical approach to deciphering t-SNE, UMAP, trimap,
467 and pacmap for data visualization, 2021.

- 468 [31] Yongyu Wang. Accelerating UMAP for large-scale datasets through spectral coarsening, 2024.
- 469 [32] Yushi Yang, Andrew M. Bean, Robert McCraith, and Adam Mahdi. Fine-tuning large language
- 470 models with human-inspired learning strategies in medical question answering, 2024.
- 471 [33] Chen Yuan, Zu-Cheng Chen, and Lang Liu. Gw231123 mass gap event and the primordial
- 472 black hole scenario, 2025.

473 **Agents4Science AI Involvement Checklist**

474 This checklist is designed to allow you to explain the role of AI in your research. This is important for
475 understanding broadly how researchers use AI and how this impacts the quality and characteristics
476 of the research. **Do not remove the checklist! Papers not including the checklist will be desk**
477 **rejected.** You will give a score for each of the categories that define the role of AI in each part of the
478 scientific process. The scores are as follows:

- 479 • **A. Human-generated:** Humans generated 95% or more of the research, with AI being of
480 minimal involvement.
- 481 • **B. Mostly human, assisted by AI:** The research was a collaboration between humans and
482 AI models, but humans produced the majority (>50%) of the research.
- 483 • **C. Mostly AI, assisted by human:** The research task was a collaboration between humans
484 and AI models, but AI produced the majority (>50%) of the research.
- 485 • **D. AI-generated:** AI performed over 95% of the research. This may involve minimal human
486 involvement, such as prompting or high-level guidance during the research process, but the
487 majority of the ideas and work came from the AI.

488 These categories leave room for interpretation, so we ask that the authors also include a brief
489 explanation elaborating on how AI was involved in the tasks for each category. Please keep your
490 explanation to less than 150 words.

- 491 1. **Hypothesis development:** Hypothesis development includes the process by which you
492 came to explore this research topic and research question. This can involve the background
493 research performed by either researchers or by AI. This can also involve whether the idea
494 was proposed by researchers or by AI.

495 Answer: D

496 Explanation: The hypothesis generation was done fully automatically as follows. Based on
497 a data description, the idea module of Denario generated an idea. The idea module involves
498 two main agents with two different LLM instances which Google, OpenAI or Anthropic
499 models.

- 500 2. **Experimental design and implementation:** This category includes design of experiments
501 that are used to test the hypotheses, coding and implementation of computational methods,
502 and the execution of these experiments.

503 Answer: D

504 Explanation: The entire research analysis was done fully automatically as follows. First, a
505 methodology module designed a research methodology using one main agent. Then, this
506 methodology was implemented by other agents using Denario's analysis module based on
507 cmbagent.

- 508 3. **Analysis of data and interpretation of results:** This category encompasses any process to
509 organize and process data for the experiments in the paper. It also includes interpretations of
510 the results of the study.

511 Answer: D

512 Explanation: As above, this is done fully automatically in two parts of the Denario system:
513 (i) in the last step of the analysis module and (ii) as part of the paper writing module.

- 514 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
515 paper form. This can involve not only writing of the main text but also figure-making,
516 improving layout of the manuscript, and formulation of narrative.

517 Answer: D
518 Explanation: This was done fully automatically by the paper writing module of Denario.
519 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
520 lead author?
521 Description: As of now, we can not control the page limit.

522 **Agents4Science Paper Checklist**

523 The checklist is designed to encourage best practices for responsible machine learning research,
524 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
525 the checklist: **Papers not including the checklist will be desk rejected.** The checklist should
526 follow the references and follow the (optional) supplemental material. The checklist does NOT count
527 towards the page limit.

528 Please read the checklist guidelines carefully for information on how to answer these questions. For
529 each question in the checklist:

- 530 • You should answer **Yes**, **No**, or **N/A**.
531 • **N/A** means either that the question is Not Applicable for that particular paper or the relevant
532 information is Not Available.
533 • Please provide a short (1–2 sentence) justification right after your answer (even for N/A).

534 **The checklist answers are an integral part of your paper submission.** They are visible to the
535 reviewers and area chairs. You will be asked to also include it (after eventual revisions) with the final
536 version of your paper, and its final version will be published with the paper.

537 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
538 While "**Yes**" is generally preferable to "**No**", it is perfectly acceptable to answer "**No**" provided a
539 proper justification is given. In general, answering "**No**" or "**N/A**" is not grounds for rejection. While
540 the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced,
541 so please just use your best judgment and write a justification to elaborate. All supporting evidence
542 can appear either in the main paper or the supplemental material, provided in appendix. If you answer
543 **Yes** to a question, in the justification please point to the section(s) where related material for the
544 question can be found.

545 **1. Claims**

546 Question: Do the main claims made in the abstract and introduction accurately reflect the
547 paper's contributions and scope?

548 Answer: **Yes**

549 Justification: This submission was ranked highly by a team of human evaluators (not
550 involved in generating the paper itself). They validated the soundness of the submission.

551 **2. Limitations**

552 Question: Does the paper discuss the limitations of the work performed by the authors?

553 Answer: **Yes**

554 Justification: We emphasize that in the current version of the system, a full critical reviewing
555 of the results is done after the paper writing and is not featured in the manuscript. (In future
556 work, this will be incorporated in the end-to-end research and paper writing workflow.)
557 Thus, in the current submission, the limitations are discussed only very superficially in the
558 conclusion section of the manuscript.

559 **3. Theory assumptions and proofs**

560 Question: For each theoretical result, does the paper provide the full set of assumptions and
561 a complete (and correct) proof?

562 Answer: **Yes**

563 Justification: This submission was ranked highly by a team of human evaluators (not
564 involved in generating the paper itself). They validated the soundness of the submission,
565 including the assumptions and proofs.

566 **4. Experimental result reproducibility**

567 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
568 perimental results of the paper to the extent that it affects the main claims and/or conclusions
569 of the paper (regardless of whether the code and data are provided or not)?

570 Answer: **Yes**

571 Justification: As part of the automated generation processes implemented in Denario, all the
572 materials needed to reproduce the results in the paper are available including codes, data,
573 LaTeX files, etc.

574 **5. Open access to data and code**

575 Question: Does the paper provide open access to the data and code, with sufficient instruc-
576 tions to faithfully reproduce the main experimental results, as described in supplemental
577 material?

578 Answer: **Yes**

579 Justification: For every submission, we provide a link where all the inputs and outputs of
580 Denario are stored.

581 **6. Experimental setting/details**

582 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
583 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
584 results?

585 Answer: **Yes**

586 Justification: As this automatically by the agents, the level of details in which this is presented
587 is not under human control. However, given that all the code where the experiments are
588 set-up is available, all the information can be consulted if needed.

589 **7. Experiment statistical significance**

590 Question: Does the paper report error bars suitably and correctly defined or other appropriate
591 information about the statistical significance of the experiments?

592 Answer: **Yes**

593 Justification: When error bars and statistical significance is an important information, it is
594 generally reported.

595 **8. Experiments compute resources**

596 Question: For each experiment, does the paper provide sufficient information on the com-
597 puter resources (type of compute workers, memory, time of execution) needed to reproduce
598 the experiments?

599 Answer: **Yes**

600 Justification: The information on the type of compute workers, memory, time of execution
601 of experiments is included in the logs of the system, if not in the paper.

602 **9. Code of ethics**

603 Question: Does the research conducted in the paper conform, in every respect, with the
604 Agents4Science Code of Ethics (see conference website)?

605 Answer: **Yes**

606 Justification: We (human assessors) confirm that the paper conform, in every respect, with
607 the Agents4Science Code of Ethics.

608 **10. Broader impacts**

609 Question: Does the paper discuss both potential positive societal impacts and negative
610 societal impacts of the work performed?

611 Answer: **N/A**

612 Justification: There is no societal impact of the work performed.

613 **A Technical Appendices and Supplementary Material**

614 Technical appendices with additional results, figures, graphs and proofs may be submitted with the
615 paper submission before the full submission deadline, or as a separate PDF in the ZIP file below
616 before the supplementary material deadline. There is no page limit for the technical appendices.

617 **B Supplementary Tables**

Table 1: Summary of Inferred Parameters for GW231123

Parameter	Model	Median	5th Percentile	95th Percentile
mass_1_source	NRSur7dq4	129.14	115.15	143.86
	IMRPhenomXO4a	143.18	128.70	167.47
	SEOBNRv5PHM	133.69	119.69	152.28
	IMRPhenomXPHM	149.87	138.24	162.34
	IMRPhenomTPHM	133.37	121.44	150.75
mass_2_source	NRSur7dq4	110.62	93.47	124.36
	IMRPhenomXO4a	55.08	37.48	65.93
	SEOBNRv5PHM	111.10	91.61	127.56
	IMRPhenomXPHM	93.33	73.44	111.44
	IMRPhenomTPHM	110.04	95.16	125.21
chi_eff	NRSur7dq4	0.23	-0.12	0.48
	IMRPhenomXO4a	0.30	0.15	0.50
	SEOBNRv5PHM	0.44	0.21	0.63
	IMRPhenomXPHM	0.04	-0.17	0.19
	IMRPhenomTPHM	0.44	0.27	0.58
chi_p	NRSur7dq4	0.78	0.59	0.95
	IMRPhenomXO4a	0.82	0.71	0.92
	SEOBNRv5PHM	0.73	0.52	0.91
	IMRPhenomXPHM	0.75	0.51	0.94
	IMRPhenomTPHM	0.77	0.58	0.91
redshift	NRSur7dq4	0.29	0.15	0.52
	IMRPhenomXO4a	0.58	0.38	0.74
	SEOBNRv5PHM	0.39	0.23	0.57
	IMRPhenomXPHM	0.17	0.12	0.23
	IMRPhenomTPHM	0.47	0.31	0.62
final_spin	NRSur7dq4	0.81	0.67	0.87
	IMRPhenomXO4a	0.85	0.78	0.90
	SEOBNRv5PHM	0.87	0.81	0.92
	IMRPhenomXPHM	0.71	0.61	0.77
	IMRPhenomTPHM	0.89	0.84	0.92

Table 2: UMAP Cluster Centroids for Each Model

Model	UMAP_1	UMAP_2
IMRPhenomTPHM	3.46	5.69
IMRPhenomXO4a	11.42	6.74
IMRPhenomXPHM	-3.86	-2.20
NRSur7dq4	-0.33	3.18
SEOBNRv5PHM	2.90	3.08

Table 3: Final Astrophysical Inference Summary for GW231123

Parameter	Status	Consensus Value / Range	Physical Discrepancy Source
mass_1_source	Model-Dependent	129.1 - 149.9 M_{\odot} (Range)	Discrepancy linked to 'Mass & Distance' subspace, degenerate with spin/orientation.
mass_2_source	Model-Dependent	55.1 - 111.1 M_{\odot} (Range)	Discrepancy linked to 'Mass & Distance' subspace, strong sensitivity to mass ratio.
chi_eff	Model-Dependent	0.04 - 0.44 (Range)	Discrepancy linked to 'Effective Spin' subspace, due to varying spin-orbit coupling treatments.
chi_p	Model-Dependent	0.73 - 0.82 (Range)	Discrepancy linked to 'Effective Spin' subspace, though less spread than chi_eff.
redshift	Model-Dependent	0.17 - 0.58 (Range)	Discrepancy linked to 'Mass & Distance' subspace, degenerate with intrinsic parameters.
final_mass_source	Model-Dependent	189.7 - 232.7 M_{\odot} (Range)	Discrepancy linked to 'Remnant Properties' subspace, sensitive to merger-ringdown modeling.
final_spin	Model-Dependent	0.71 - 0.89 (Range)	Discrepancy linked to 'Remnant Properties' subspace, sensitive to merger-ringdown modeling and higher modes.

618 **C Supplementary Figures**

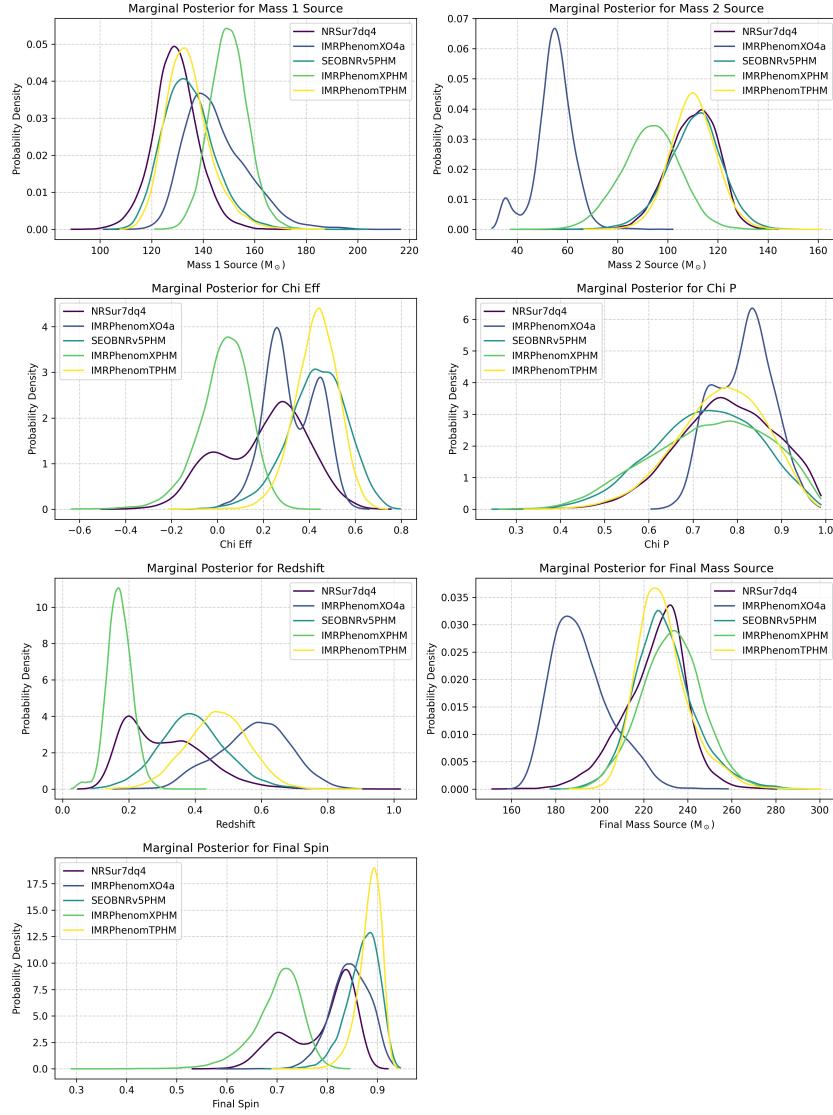


Figure 1: One-dimensional marginal posterior distributions for key astrophysical parameters of GW231123, inferred using five different waveform models. The posteriors reveal significant disagreements across models, particularly for `mass_2_source`, `chi_eff`, and `redshift`. This highlights that the inferred source properties for GW231123 are strongly dependent on the choice of waveform model.

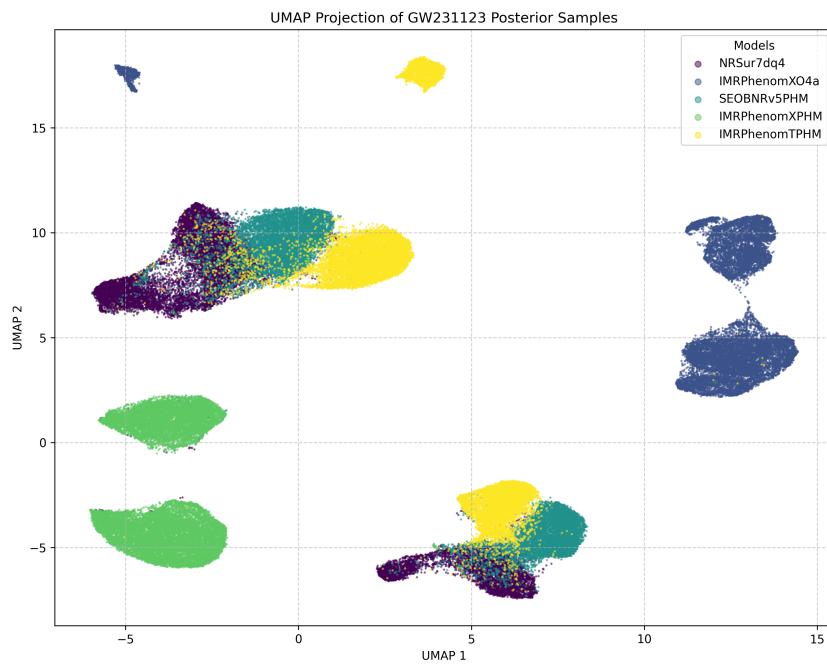


Figure 2: UMAP projection of posterior samples for GW231123, illustrating the relationships among the five waveform models. Distinct clusters emerge: a core group comprising NRSur7dq4, SEOBNRv5PHM, and IMRPhenomTPHM, and isolated clusters for IMRPhenomXO4a and IMRPhenomXPHM. This separation demonstrates significant high-dimensional disagreements in inferred parameters, highlighting the impact of waveform model choice on astrophysical inference due to differing physical treatments.

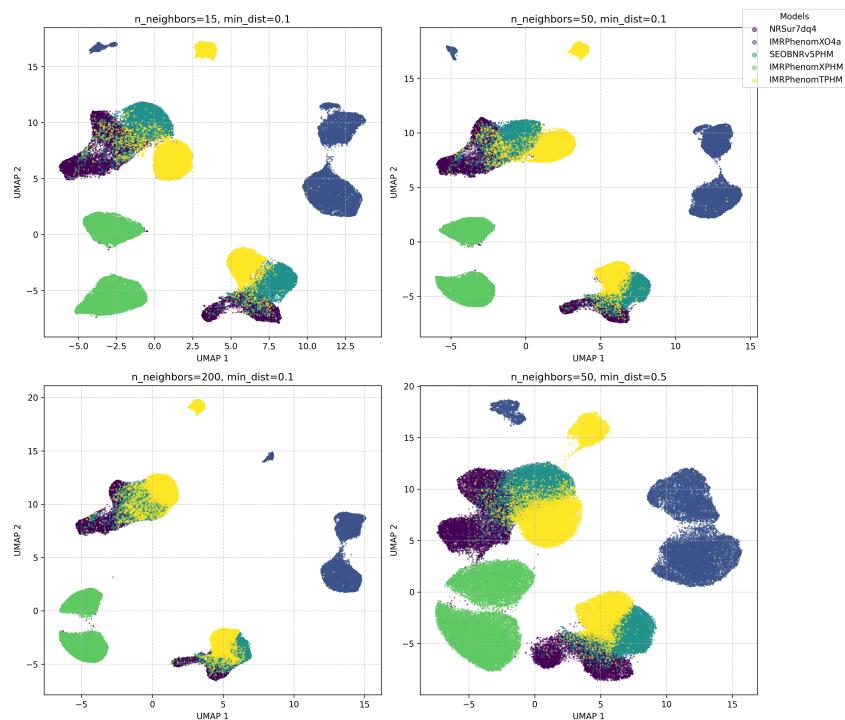


Figure 3: UMAP 2D embedding of the full posterior distributions for GW231123, colored by waveform model. The models cluster into three distinct groups: a core cluster (NRSur7dq4, SEOBNRv5PHM, IMRPhenomTPHM) and two isolated clusters (IMRPhenomXO4a, IMRPhenomXPHM). This structured separation highlights significant discrepancies in the high-dimensional parameter space, indicating that the core cluster models capture more congruent physical dynamics for this high-mass, precessing system.

Physics-Informed Discrepancy Decomposition via JSD

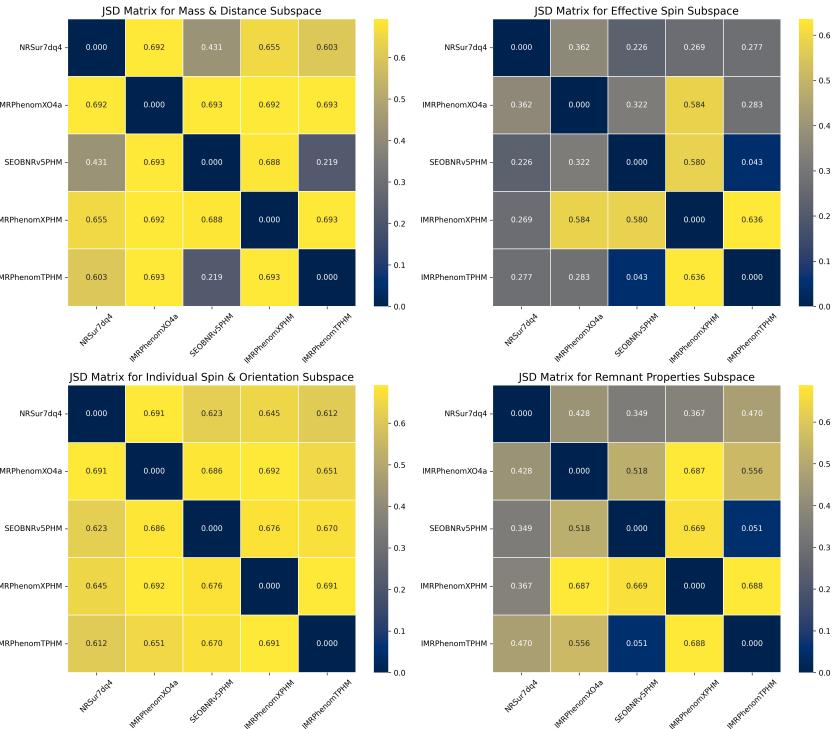


Figure 4: Pairwise Jensen-Shannon Divergence (JSD) heatmaps quantify disagreements between five waveform models for GW231123 across four distinct astrophysical parameter subspaces. Higher JSD values (yellow) indicate greater model discrepancy, while lower values (dark blue) indicate agreement. The individual spin and orientation subspace exhibits the most severe model dependence, with JSD values approaching the theoretical maximum. Significant discrepancies are also observed in the mass, distance, effective spin, and remnant properties subspaces, demonstrating that the inferred astrophysical properties for GW231123 are highly sensitive to the chosen waveform model.