# Neural surrogates for designing gravitational wave detectors

Carlos Ruiz-Gonzalez[1,2], Sören Arlt[1,3], Sebastian Lehner[2], Arturs Berzins[2], Yehonathan Drori[4,5], Rana X Adhikari[4], Johannes Brandstetter[2,6], and Mario Krenn[1,3]

[1]Max Planck Institute for the Science of Light, Erlangen, Germany

[2]ELLIS Unit Linz, Institute for Machine Learning, JKU Linz, Austria

[3]Machine Learning in Science Cluster, Department of Computer Science, Faculty of Science, University of Tuebingen, Germany

[4]LIGO Laboratory, California Institute of Technology, Pasadena, USA

[5]School of Physics and Astronomy, Tel Aviv University, Israel

[6]Emmi AI, Linz, Austria

**Physics simulators are essential in science and engineering, enabling the analysis, control, and design of complex systems. In experimental sciences, they are increasingly used to automate experimental design, often via combinatorial search and optimization. However, as the setups grow more complex, the computational cost of traditional, CPU-based simulators becomes a major limitation. Here, we show how neural surrogate models can significantly reduce reliance on such slow simulators while preserving accuracy. Taking the design of interferometric gravitational wave detectors as a representative example, we train a neural network to surrogate the gravitational wave physics simulator Finesse, which was developed by the LIGO community. Despite that small changes in physical parameters can change the output by orders of magnitudes, the model rapidly predicts the quality and feasibility of candidate designs, allowing an efficient exploration of large design spaces. Our algorithm loops between training the surrogate, inverse designing new experiments, and verifying their properties with the slow simulator for further training. Assisted by auto-differentiation and GPU parallelism, our method proposes high-quality experiments much faster than direct optimization. Solutions that our algorithm finds within hours outperform designs that take five days for the optimizer to reach. Though shown in the context of gravitational wave detectors, our framework** is broadly applicable to other domains where simulator bottlenecks hinder optimization and discovery.

Carlos Ruiz-Gonzalez: cruizgo@proton.me

Mario Krenn: mario.krenn@uni-tuebingen.de

## 1 Introduction

Computer simulations are ubiquitous in all branches of science and engineering, providing predictions and a deeper understanding of complex interacting systems [1–3]. In experimental sciences, they are also used to optimize, control, and even design experimental setups. This goes beyond tuning a few parameters; novel algorithms have (partially) automated the conception of experiments, leading to new and often counter-intuitive designs [4]. In physics alone, there are successful examples in quantum physics [5–11], particle physics [12–14], nano-optics [15–17], microscopy [18, 19] and gravitational wave physics [20], to name a few.

The inverse design of experiments is extremely computationally intensive [21], especially when combinatorial search is required. Therefore, the simulator speed is paramount to navigate the vast space of experimental designs. Fortunately, these high computational demands are mitigated by hardware improvements, such as GPUs [22, 23], and high-performance tools, like auto-differentiation or just-in-time compilation. These advances can also benefit slower CPU-based software, which can be surrogated with neural networks [24–31].

In this work, we leverage fast neural surrogate models to design interferometric gravitational wave detectors (GWDs). These sensitive experiments, able to detect distant cosmological events, can be simulated with the CPU-based

software Finesse [32]. This open-source simulator was developed by the LIGO collaboration, the first to detect gravitational waves in 2015 [33], and is the de facto standard to GWD design. In recent work [20], Finesse was used to automate the search of new interferometric detectors by optimizing a highly parametrized experimental ansatz [20]. However, as the Finesse simulator is neither differentiable nor GPU accelerated, gradient-based optimization becomes increasingly inefficient when the system grows.

To accelerate the design process, we present SPROUT, short of *Surrogate Predictions for Reiterative Optimization with Update Training.* Our multi-step approach trains a surrogate model to predict the quality and feasibility of different GWDs, encoded as the parameters of a UIFO. First, we train the model with randomly parametrized designs, to then generate new data that, once verified, is used for further training. Iterating the design and training (see Figure 2), we produce successive generations of models that efficiently explore the design space and improve the solutions. This iterative optimization process, also known as online optimization or active learning [34, 35], is a successful strategy when simulations (or experiments) are expensive [36–39].

Our approach does not replace the original simulator, as we still use it to verify the designs generated by the faster GPU-based differentiable surrogates. However, by rapidly exploring the space of possible designs, even if it is an approximation, we reduce the overall usage of the CPU software and still obtain high-quality designs. Further, since the verifications of the GPU-based optimization are independent of each other, we can parallelize the simulator calls arbitrarily, beyond what is possible in a sequential gradient-based optimization. Finally, when the size of the GWD increases, the growing number of tunable parameters and the longer simulation times make the numerical gradient approximations prohibitively slow. Neural surrogates excel in such a scenario, as backpropagation and the vast parallelization capabilities of GPUs allow us to optimize thousands of experiments at breakneck speed. We obtain high-quality designs in a fraction of the time, and, remarkably, even when we train the models with randomly parametrized designs, they find solutions far beyond anything we trained them with (see Figure 2b).
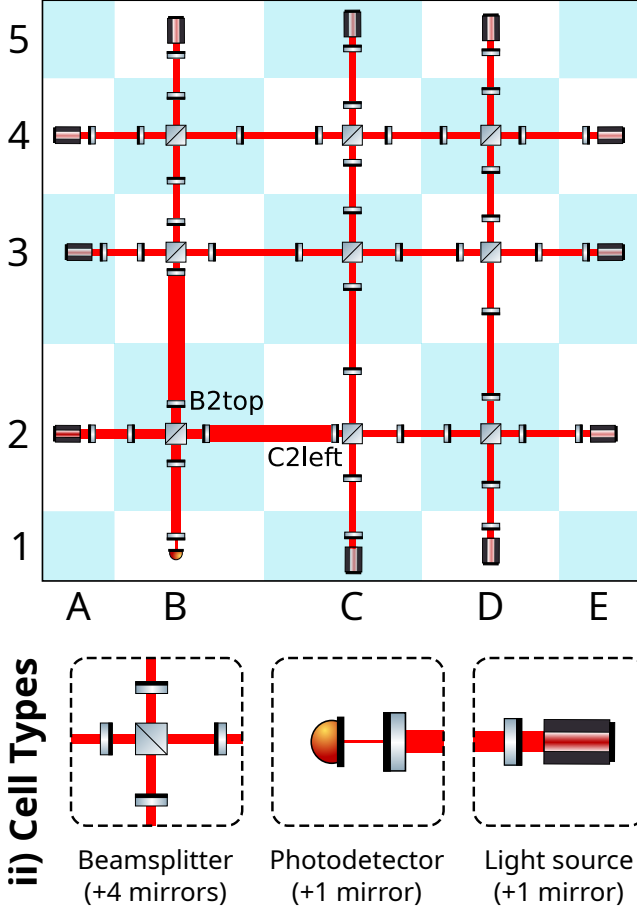
## 1.1 Related works

**Surrogates for inverse design.** When the computational cost of traditional simulators becomes unbearable, (neural) surrogate models are efficient, often differentiable, alternatives [40]. They are widely used for inverse design tasks in multiple fields, including engineering [41–43], chemistry [44–48], or physics [30, 31, 49–51], being particularly popular in photonics [24–27, 52]. Most of these surrogates predict one or a few values that summarize the designs' quality. In contrast, we predict an array of properties that indicate the quality as well as the hardware requirements to build the design. These predicted properties are incredibly responsive to the input parameters, and can fluctuate up to 10 orders of magnitude when tuning a single experimental parameter (see Figure 1). We did not find such sensitive systems in the surrogate literature and took inspiration from other domains [53] to face the challenge (details on Section 2).

**Neural networks in physics.** Recent advances in machine learning are especially aimed at the physical sciences. Physics-informed neural networks (PINN) [54, 55], or neural operators [56, 57], are notable examples, widely used as PDE surrogates in fluid dynamics [58, 59], among other domains [60, 61]. However, these novel approaches are not suited to our task: predicting the experiment's sensitivity and the light powers applied on the elements, based on the design parameters. For the PINNs, we lack simple physical laws to constraint the outputs. On the other hand, neural operators surrogate functions which can be evaluated at any discretization of a continuous input space. Here we have a finite set of input properties, from a finite number of optical parameters, and a fixed number of output properties. Thus, these tools offer us no advantage.

**Pool-based active learning.** Supervised learning requires costly labeled data. Thus, active learning strategies must carefully select which data to label, from a pool of unlabeled data points [34]. There exist theoretical results about the strengths and limitations of several approaches [62–65], but their assumptions rarely hold when training deep neural networks. Therefore, the field relies mostly on heuristics and empirical evidence for particular applications [66–69].

## i) UIFO gravitational wave detector



## ii) Cell Types



Beamsplitter (+4 mirrors)  Photodetector (+1 mirror)  Light source (+1 mirror)
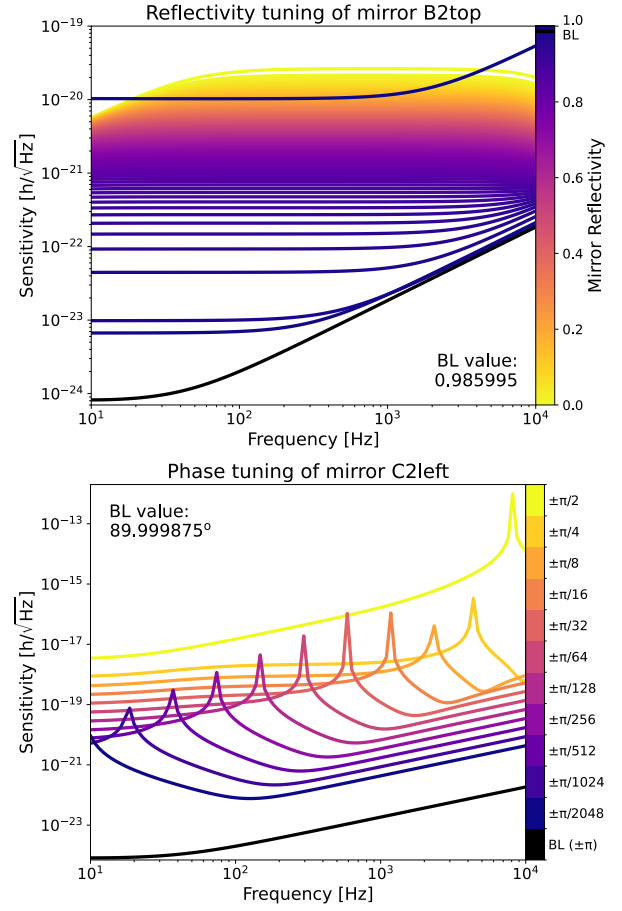
## iii) Highly responsive solutions



Figure 1: **A highly expressive ansatz: quasi-Universal InterFerOmeter (UIFO).** A large variety of gravitational wave detectors can be built by aligning optical elements in an irregular rectangular grid. This expressive parametrized template is referred to as quasi-Universal InterFerOmeter, or UIFO. **i)** The 3x3 UIFO, with 169 parameters, has 9 beamsplitter cells surrounded by light source cells and a single photodetector. By tuning the parameters of the grid's components, using only a small subset of components, we can recreate a simplified version of the GWD used by the LIGO collaboration. **ii)** In the depicted design and throughout this paper, we use beamsplitters, mirrors, photodetectors, and lasers. **iii)** The quality of a GWD is given by the smallest spatial deformation that it can detect for each gravitational wave frequency, represented by a sensitivity curve. As shown in the top plot, these curves can fluctuate in several orders of magnitude in response to a single parameter, like the reflectivity of the top mirror in cell B2. The change is even more drastic when tinkering with the left mirror in cell C2.

Bayesian methods are also popular for small systems [70–73], but they do not scale well. In this work, we verify the designs obtained after a noisy optimization process of the surrogate simulator, pushing the model to generalize beyond its training dataset. We repeat this procedure for 5 rounds (see Figure 2a). Our strategy resembles previous work on active learning [74, 75] and on adversarial training [76], where gradient-based optimization generates adversarial examples that are mislabeled by neural [77, 78]. Retraining the networks with such examples increases their robustness and improves their predictions in high-quality design spaces.

## 2 Methods

This section describes how to design GWD using the *quasi-Universal InterFerOmeter* (UIFO), a highly overparametrized experimental setup introduced in previous work [20], which is used as a highly expressive ansatz. We also outline the surrogate models' architecture and training strategy, the building blocks of the SPROUT algorithm.

### 2.1 Gravitational wave detectors

Gravitational waves are perturbations of the spacetime caused by the motion of masses. They

were predicted by Einstein's general relativity in 1916 [79] and directly observed for the first time in 2015 by the LIGO collaboration [33]. Feeble as they are, we can only detect the waves originating from massively energetic astrophysical phenomena, such as supernovae or the merging of black holes. To do so, experimentalists use some of the most sensitive devices ever created by humans – gravitational wave detectors – massive facilities whose components require extremely precise tuning. Yet, as shown in previous work [20], the space of GWD experiments is mostly unexplored. Many possible designs, including the LIGO detector, can be described within the UIFO framework, a rectangular lattice of parametrized optical elements. Using an expressive ansatz is a powerful tool to design new experiments, also in other fields [9, 19, 80], as it relaxes the optimization problem from (partially) discrete to purely continuous. The price we pay is to optimize larger systems of tunable elements, likely redundant, which contain previously known baselines and new designs able to surpass them.

As shown in Figure 1, for a given choice of UIFO parameters, the detector will have a certain strain sensitivity for each possible frequency of a gravitational wave, which we refer to as *sensitivity curve*. The lower the sensitivity, the smaller space fluctuations it can detect. Additionally, it is necessary to limit the optical power in each component, otherwise their performance degrades. The same applies to the photodetector, which is orders of magnitude more sensitive. These and other constraints are summarized in Table 1.

The sensitivity and power values define the quality and viability of the designs and can be computed with Finesse, the CPU-based simulator we want to surrogate. As shown in Figure 1iii), the quality of the interferometer can fluctuate by multiple orders of magnitude when tuning even a single parameter. Furthermore, some of the plotted sensitivity curves could not be implemented with current hardware, as the light power at several components surpass the forementioned limitations.

In this project, we focus on the continuous optimization of the parametrized ansatz UIFO and how the benefit of the surrogates increases with the size of the GWD. Therefore, to facilitate the comparison between designs of different sizes, we fix all the discrete design variables: the grid size, the beam splitters' orientations, and the location of the photodetector. The reduced size of our setups also facilitates the experiments. For more practical applications, we would emulate larger designs with different configurations that include even more sensitive designs than the LIGO detector. For further details on the UIFO parametrization, see Appendix A.

| Parameter | Min | Max |
|---|---|---|
| Optical path length | 1m | 4km |
| Optical loss per element | 5ppm | |
| Optical transmission | 15ppm | |
| Reflected optical power | | 3.5MW |
| Transmitted optical power | | 2kW |
| Power in photodetector | | 10mW |

Table 1: **Ranges of the physical parameters of the UIFO designs.** The loss/transmission of the optical elements can be bounded before running the simulations. Same for the optical path length, which ensures that the design remains within the 4km that LIGO design measures. The powers at the optical elements, computed by the simulator, are enforced via penalties. The lower bound for the optical loss, 5ppm, is a rough estimate of what is possible with current technology. We assume the lowest possible loss for all components, as larger values would only worsen the performance. Further details in Appendix A.
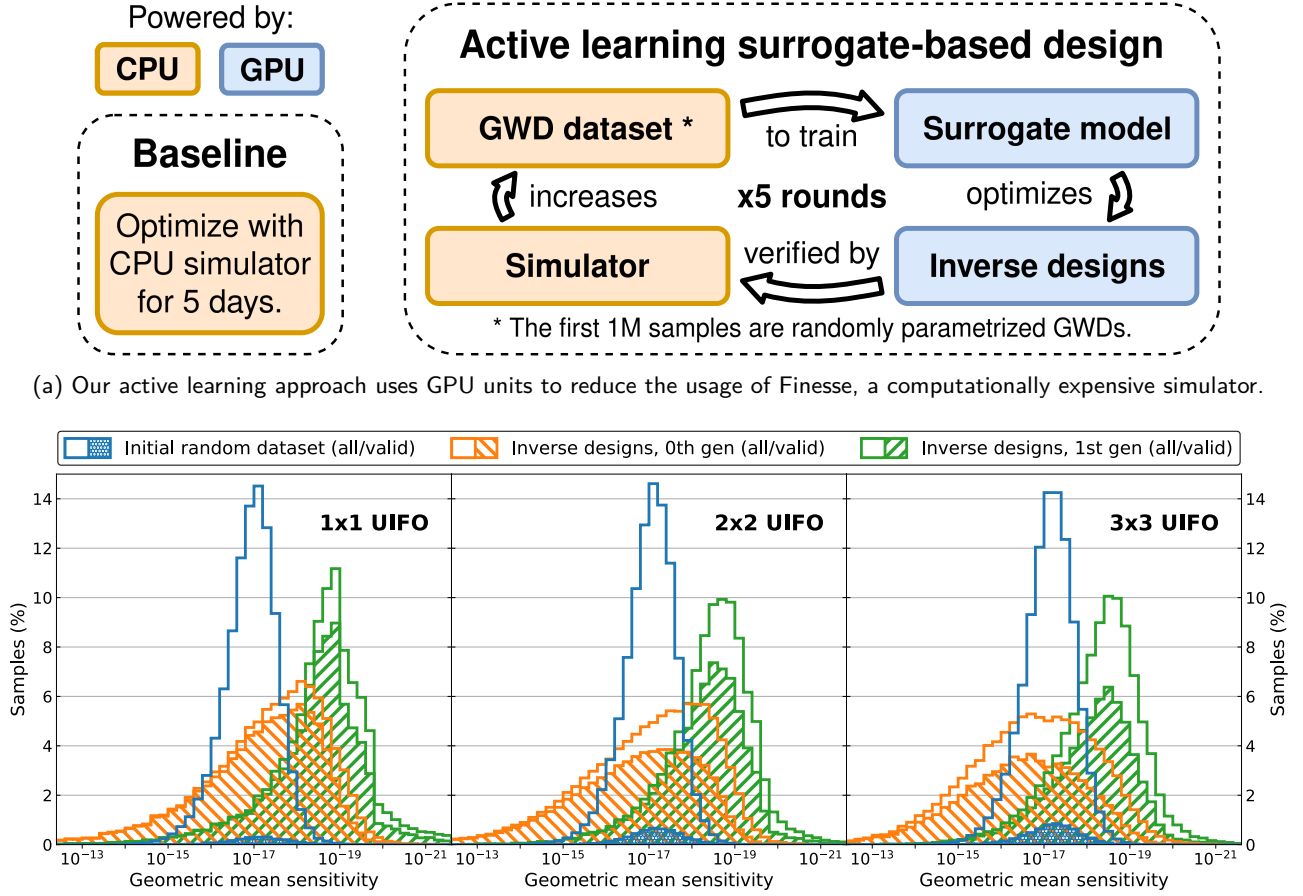
## 2.2 Design of experiments

Whether we simulate the UIFO with Finesse or rely on a neural surrogate, the strategy to design novel GWD is the same. Having fixed the discrete design choices, we optimize the continuous parameters of the UIFO, our template for GWD, to minimize the geometric mean of the strain sensitivity.
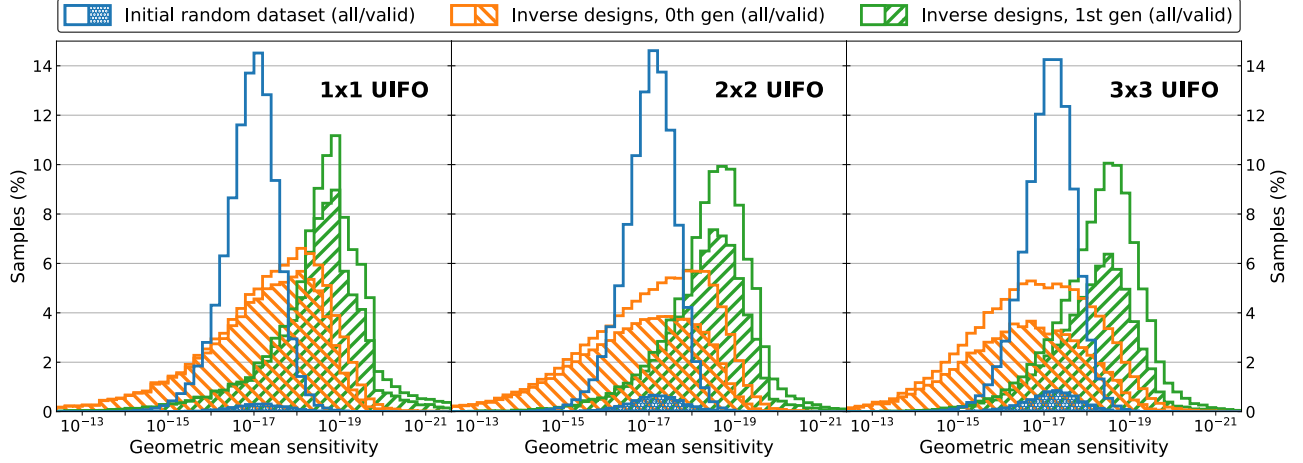
Furthermore, to guarantee the feasibility of the designs, they must meet the hardware constraints described in Table 1, which can be divided into two types. First, we have the maximum powers that the detectors, mirrors, and beam splitters can withstand. These are output constraints that are calculated by the simulator. We enforce them by adding penalties to the loss, using the differentiable function

$$\text{Penalty} \equiv (\alpha + \beta(x - x_{\text{th}}))\sigma(\gamma(x - x_{\text{th}})), \quad (1)$$

which approximates a step function with a step

(a) Our active learning approach uses GPU units to reduce the usage of Finesse, a computationally expensive simulator.



(b) The designs' quality (i.e. low sensitivity) surpasses the model's training data. Only valid samples meet hardware constrains.

| Experiment size | 1x1 (29 parameters) | | 2x2 (85 parameters) | | 3x3 (169 parameters) | |
|---|---|---|---|---|---|---|
| Design strategy | Simulator | Surrogate | Simulator | Surrogate | Simulator | Surrogate |
| Processor unit | CPU only | CPU+GPU | CPU only | CPU+GPU | CPU only | CPU+GPU |
| Time until 1e-20 | 1h10' | 14'+40' | 6h05' | 41'+1h09' | 32h15' | 1h23'+2h01' |
| Simulation calls to reach 1e-20 | 5.5M | 2.2M | 9.6M | 2.2M | 25.6M | 2.2M |

(c) For the same quality, the surrogate-based design (+ random sampling) takes much less time than only using the simulator.

Figure 2: **Iterative inverse design with surrogate models.**

size $\alpha$ and a final slope of $\beta$, for $x > x_{\text{th}}$. The stiffness of the step is indicated with $\gamma$, $\sigma$ is a sigmoid function, $x$ is the (log) optical power applied to a given optical element, and $x_{\text{th}}$ is the maximum power that the element can tolerate. The values of $x_{\text{th}}$ depend on the optical element and whether the power is reflected or transmitted (see Table 1). We set the values of $\{\alpha, \beta, \gamma\}$ to $\{10, 1, 1000\}$.

In contrast, the constraints on element properties are enforced by restricting the input parameter space, for example, by limiting the reflectivity and losses of the mirrors and beam splitters. To guarantee spatial consistency, the distances

between rows and columns of beamsplitters are described as fractions of a total larger distance, which must remain below 4 km. Finally, the positions of the different mirrors are also described by fractions (see Appendix A).

The final loss is the mean of the log sensitivity for $N$ frequency values $f_i$ (i.e. the log of the geometric mean) plus the described penalties,

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} \log(\text{S}(f_i)) + \text{Penalties}. \qquad (2)$$

The set of frequencies are 101 logarithmically spaced values between 1Hz and 10kHz.

To minimize the loss, we use the Adam optimizer [81] with learning rate of $\lambda = 0.1$. Moreover, we add annealed Gaussian noise, $\mathcal{N}(0, 2\epsilon^t \lambda)$, to the input update, a well-known strategy to optimize highly non-convex functions [82–84]. With $\epsilon < 1$, the noise decreases exponentially at every step $t$.

## 2.3 Surrogate models

In order to find new GWD designs, our neural networks must be able to simulate them, at least approximately. We want our surrogates to predict, based on the input parameters of the UIFO, the final sensitivity curve of the designs (their quality) and the light powers that are applied to each optical component (their feasibility). To process the grid-like structure of the UIFO template, we divide the interferometer into overlapping subgraphs, each with a different set of parametrized elements, and turn them into patches, which are then processed by a transformer encoder and mapped into the output space. Our architecture is heavily inspired by models such as Vision Transformers [85], MLP-Mixers [86], and other works that structure the input data into patches/tokens, [87]. It has three main parts:

1. the patch extraction from the GWD and their embedding into tokens,

2. the processing of the tokens with a standard transformer encoder of 8 layers, and

3. the decoding of the processed tokens into the final output of sensitivities and powers.

**Turning GWD into patches.** To transform a GWD into a set of patches, we treat the UIFO grid as an attributed graph, i.e. a graph whose nodes and edges have several weights. The nodes are the beam splitters, the light sources, and the detectors, each with a different set of associated parameters. Between the nodes, we always find two mirrors, whose properties and relative position in the optical paths are the attributes of the edges that connect the nodes. Grouping the parameters of an edge and the two adjacent nodes, we define a patch. Since the nodes have different numbers of attributes, we pad the patches with zeroes up to length 12. As shown in Figure 3, the patches overlap in the beam splitters, as their respective nodes have degree 4.

To the padded patches, we add a positional embedding (PE) and apply the trainable Fourier feature map, that is, a linear layer followed by a sine activation [53]. This mapping helps to overcome the spectral bias of neural networks towards low-frequency functions [88, 89], which can limit the capacity of our models to predict the properties of high-quality outliers. The projection of the Fourier map increases the embedding dimension of the patches to 192. It is followed by a linear transformation.

**Transformer architecture.** After adding a positional embedding to the tokens, we pass them to the transformer encoder [90]. We use eight standard transformer layers with the pre-layer normalization convention [91], and a multilayer perceptron (MLP) with hidden dimension 384, twice the size of the tokens. We chose the transformer architecture, as it has shown great success processing multiple types of data for a wide variety of tasks [92–100].

**Towards the final output.** After the transformer encoder, we apply a layer norm and reduce the tokens' size with a linear projection, followed by a GELU. Concatenating the tokens, we obtain a final vector that, after a norm layer and an MLP, takes the shape of the output. We enforce normalization by applying a sigmoid to the output: the array of sensitivities and the light powers.

A summary of the architecture is shown in Figure 3. The implementation details are described in Appendix C, and all the code can be found in the repository[1].

The training datasets consist of pairs of input parameters and output sensitivities and powers. The input parameters are linearly normalized between 0 and 1, adjusting the rescaling to the different types of physical parameters. On the other hand, since the values of sensitivities and optical powers range across several orders of magnitude, we take their logarithm value for training. For further details on data preprocessing, see Appendix B.

## 2.4 Active learning strategy

To train the models, we use the Adam optimizer with weight decay and a plateau learning rate schedule. In the first training phase, the model
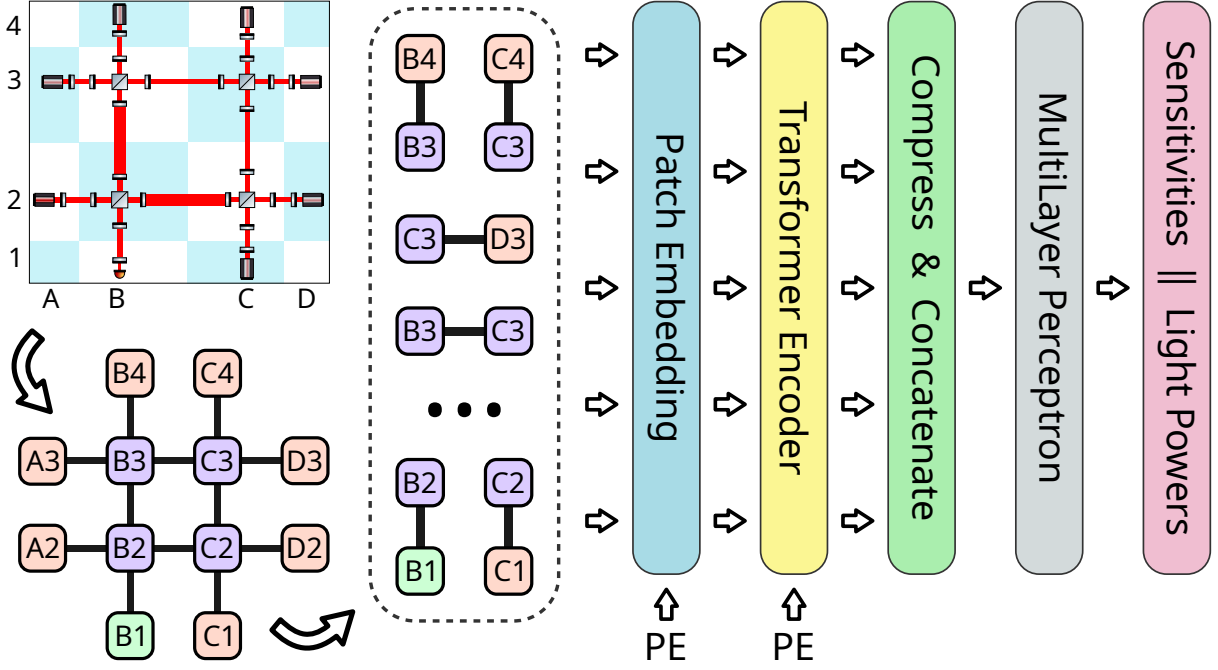
Figure 3: **Overview of the surrogate model's architecture.** After extracting the patches from the grid-based GWD, we transform them into (larger) token representations, which are then processed by a standard transformer encoder. To turn the tokens into the final output, we compress them with a learned projection and concatenate them. The resulting array, larger than the output by design, goes through a final MLP and, after a sigmoid, is turned into the output: the values of the sensitivity curve concatenated with the light powers at each optical element.

learns from 1 million randomly parametrized designs.[2] With the trained models we run the Adam optimization process described in section 2.2, adding annealed noise with $\epsilon = 0.999$. Running 200,000 optimization trajectories for 16,384 steps on multiple GPU nodes, we store 5 designs for each trajectory[3] producing one million new training samples each round (+10% for testing). After verifying the designs – in parallel – with the CPU-simulator Finesse, we use the old and the new samples, 2 million in total, to retrain the model. The total amount of training samples increases by 1 million at every training round, which are added to the training datasets in groups of 50K to stabilize training [101]. The neural networks are initialized with the weights of the previous training phase.

It is noteworthy that the quality of the first training dataset, randomly produced, is always very low. Less than 5% are valid, as most of them apply excessive light power in at least one optical component, usually the fragile detector.

Yet, the first trained model produces valid designs much more sensitive than any sample it was trained with (see Figure 2). That happens for all system sizes. After further training, the model produces increasingly better samples by learning from its previous designs, a virtuous cycle of self-improvement (see Figure 4).

## 3 Experiments

To benchmark our surrogate-based algorithm, SPROUT, we compare it with the previous baseline: direct optimization with Finesse, the CPU-based simulator. Table 2 summarizes the most relevant feats of both approaches, while Figure 4 shows the evolution of the surrogates' designs across several iterations of (re-) training and designing, comparing them with the CPU-powered optimization process. We also detail the resources employed for both approaches.

### 3.1 CPU-based optimization baseline

The original approach we want to surpass is the gradient-based optimization with the CPU-simulator Finesse. Therefore, we optimize 256

---

[2]Additional 10% are produced for testing.

[3]We store the designs at (approximately) logarithmically spaced steps: 6, 48, 337, 2352, and 16384.
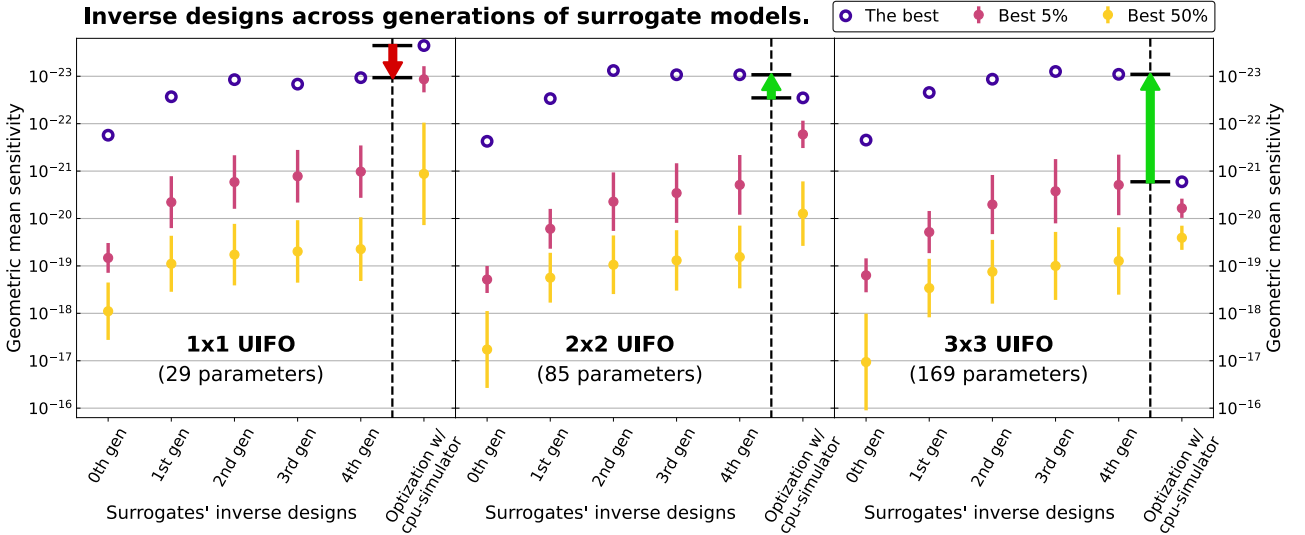
Figure 4: **Evolution of the inverse designs' quality after (re)training the surrogate model and comparison with direct optimization.** After training the surrogate model with random data we use it to produce hundreds of thousands of novel solutions via gradient-based optimization. Once the designs are verified, they are added to the previous training dataset, to finetune the model. At every generation of finetuned models the training data adds 1 million of new samples. The plotted sensitivities, which improve after every generation, come from the last step of the optimization process of the 200K inverse-designed samples. The optimizations with Finesse, detailed in Figure 7, were limited to 256 samples per grid size and ran for 5 days. The benefits of the surrogate models are specially prominent in the larger systems.

randomly parametrized UIFO samples with the Adam algorithm, adding annealed noise of $\epsilon = 0.999$, just as with the surrogates. The main difference from surrogate optimization is the need to compute the gradients with finite differences, a much slower procedure than automatic differentiation. We ran the optimization for 5 days and only the smallest UIFO with 29 parameters was able to reach more than 10 thousand optimization steps.

The final designs' sensitivities optimized with Finesse are shown in Figure 4 for the 3 UIFO sizes, together with the designs of the surrogate models. The complete evolution of the CPU-based optimization over the 5 days is shown in Figure 7.

## 3.2   The SPROUT algorithm

SPROUT overcomes the limitations of the CPU simulator by delegating part of the task to powerful and parallelizable GPU nodes. At each round, we train the model with a GPU node from *AMD Instinct MI300A APU*, and for inverse design we use 11 of such nodes. We optimize the designs in batches of 20000, for all UIFO sizes, minimizing the loss from equation 2, as described in Section 2.4.

Storing five optimized samples for every optimization trajectory, we have 1 million of new input samples for training (+10% for testing). To verify them, we run Finesse in parallel jobs across 256 CPU cores from several *Xeon Gold 6130*. That is the same hardware used for the optimization with Finesse described in Section 3.1.

To ensure the reliability of our approach, we applied SPROUT twice for each size of the UIFO. Therefore, the models are initially trained with the same random dataset but, due to the different random seed, they present minor differences in the training process that also affect the design (see Appendix D). Nonetheless, the overall quality is usually very similar. The inverse designs shown in Figure 4 are the aggregated data of both experiments, while the times in Table 2 are averaged from the two processes.

## 3.3   Results

As shown in Table 2, the advantage of the surrogate models for inverse design increases with the size of the experiment. For the largest UIFO, the designs from the surrogate model improve the sensitivities by more than an order-of-magnitude.

| UIFO grid size | 1x1 (29 parameters) | | 2x2 (85 parameters) | | 3x3 (169 parameters) | |
|---|---|---|---|---|---|---|
| Design tool | Simulator | Surrogate | Simulator | Surrogate | Simulator | Surrogate |
| Simulator calls | 560.6M | 6.6M | 189.3M | 6.6M | 95.2M | 6.6M |
| Final samples | 256 | 1.1M | 256 | 1.1M | 256 | 1.1M |
| Steps / sample | 73001 | 16384 | 8600 | 16384 | 2189 | 16384 |
| Best sensitivity | <u>2.25e-24</u> | 1.07e-23 | 2.85e-21 | <u>7.56e-24</u> | 1.68e-21 | <u>7.92e-24</u> |
| Time spent | **CPU** | **CPU+GPU** | **CPU** | **CPU+GPU** | **CPU** | **CPU+GPU** |
| Until 1e-20 | 1h10' | 14'+40' | 6h05' | 41'+1h09' | 32h15' | 1h23'+2h01' |
| Until 1e-22 | 5h09' | 28'+1h44' | 90h13' | 1h23'+2h56' | Never | 2h46'+4h37' |
| Total | 5 days | 1h24'+7h36' | 5 days | 4h10'+11h27' | 5 days | 8h18'+18h6' |

Table 2: **Time, quality, and resource comparison of the 2 design strategies for 3 sizes of UIFO.** Our surrogate-based approach drastically reduced the simulator calls to achieve the same design quality. The CPU-powered optimization ran for 5 days (fixed) in 8 units of *Xeon Gold 6130*, with 32 CPU cores each. With the same CPU resources, the verification of the 6.6M samples (5 rounds of inverse design + starting random dataset) takes only a few hours. The additional GPU time comes from the training time, which increased at every round with the size of the dataset (see Table 4), together with the inverse design, which took (approximately) 16, 40, and 85 minutes, for the 3 respective sizes. The inverse design was distributed across 11 GPU cores from multiple *AMD Instinct MI300A APU* units, each optimizing 20K samples in parallel. Details on the computational resources in Appendix D.

The increasingly better designs presented in Figure 4 show how the models are able to generalize beyond the data they are trained on. This is particularly noticeable at the first iterations (see Figure 2b).

The key behind the surrogates' performance is to utilize the GPU and the differentiability of the neural nets, being able to simulate many more experiments in a fraction of the time than required by the CPU-based simulator Finesse. To give a loose comparison, for the larger 3x3 UIFO, an optimization step takes less than 0.3 seconds on a *AMD Instinct MI300A APU* GPU core for 20,000 samples. A single evaluation with Finesse takes around 1.2 seconds on a *Xeon Gold 6130* CPU node. Needing 169 additional simulator calls to compute the gradient, every optimization step can take more than 3 minutes. Therefore, as the size of the optimized GWD grows, so do the benefits of the surrogate models.

## 4 Discussion

We introduced SPROUT, a surrogate-based approach to accelerate the inverse design of gravitational wave detectors. Our algorithm navigates the space of experimental designs by iteratively training neural surrogates, optimizing designs, and selectively verifying candidates with Finesse, the non-differentiable, CPU-based simulator that we are surrogating. While tested with a limited set of experimental tools, the surrogates exhibited outstanding generalization capabilities, as they designed experiments of much higher quality than those used for training. Given the drastic fluctuations in the GWD's properties when modifying the input parameters, the feat is even more surprising.

By exploiting differentiability and GPU-based parallelization, the benefits of our neural network approach rise with the system size, as finite-difference methods become increasingly expensive. This approach transcends our field, as it can accelerate any design task that relies on non-differentiable, CPU-based simulators. Moreover, the patch-based encoding of experimental designs is a simple yet versatile approach, suited to all types of modular setups. Finally, our active learning strategy needs little domain expertise, only clear goals and constraints for the models to pursue.

With the positive results of our algorithm, there is still room for improvement. First, our inverse designs use a small set of optical elements and a few simple topologies. Constraining the task, we can compare the performance of different methods in a well-known arena, but we are limited to rediscovery and prone to stagnation. To reach SOTA detectors, we need other components and topologies. That leads us to the second point: our neural networks could be

more flexible. We trained our neural networks on fixed discrete choices, but better models could pick different components, optimize the topology of experimental setups, or even predict the detector's sensitivity for arbitrary wave frequencies. To meet these ambitious goals, our architectures must map between different input and output spaces, and we need richer training data that encodes every possible design choice. These modifications might lead to general changes in the SPROUT algorithm, our third point. As the search space grows in dimensions and complexity, we might need more efficient training and exploration strategies. Some options include limiting the training time, intelligent sample weighting to favor specific designs' qualities, or implementing curriculum training with increasingly complex designs. Ensemble methods can foster exploration by aggregating predictions across diverse surrogate models, and diffusion-based design could be a powerful asset to explore different topologies.

To conclude, SPROUT is a promising tool for inverse design, whose principles and benefits apply beyond the search for GWD. Having showcased its potential in a constrained task, this work paves the way for further advances that lead to novel and more powerful experiments.

## Acknowledgements

## References

[1] N Metropolis, J Howlett, and G Rota. *A history of computing in the twentieth century.* Academic Press, New York, NY, 1980. DOI:10.1016/C2009-0-22029-0.

[2] Giovanni Battimelli, Giovanni Ciccotti, and Pietro Greco. *Computer Meets Theoretical Physics: The New Frontier of Molecular Simulation.* Springer International Publishing, 2020. DOI:10.1007/978-3-030-39399-1.

[3] Jack Dongarra and David Keyes. The co-evolution of computational physics and high-performance computing. *Nature Reviews Physics*, 6(10), 2024. DOI:10.1038/s42254-024-00750-z.

[4] Mario Krenn, Manuel Erhard, and Anton Zeilinger. Computer-inspired quantum experiments. *Nature Reviews Physics*, 2(11), 2020. DOI:10.1038/s42254-020-0230-4.

[5] Mario Krenn, Mehul Malik, Robert Fickler, Radek Lapkiewicz, and Anton Zeilinger. Automated search for new quantum experiments. *Physical review letters*, 2016. DOI:10.1103/PhysRevLett.116.090405.

[6] PA Knott. A search algorithm for quantum state engineering and metrology. *New Journal of Physics*, 18, 2016. DOI:10.1088/1367-2630/18/7/073033.

[7] L O'Driscoll, Rosanna Nichols, and Paul A Knott. A hybrid machine learning algorithm for designing quantum experiments. *Quantum Machine Intelligence*, 2019. DOI:10.1007/s42484-019-00003-8.

[8] Alexey A Melnikov, Pavel Sekatski, and Nicolas Sangouard. Setting up experimental bell tests with reinforcement learning. *Physical review letters*, 2020. DOI:10.1103/PhysRevLett.125.160401.

[9] Carlos Ruiz-Gonzalez, Sören Arlt, Jan Petermann, Sharareh Sayyad, Tareq Jaouni, Ebrahim Karimi, Nora Tischler, Xuemei Gu, and Mario Krenn. Digital discovery of 100 diverse quantum experiments with pytheus. *Quantum*, 7:1204, 2023. DOI:10.22331/q-2023-12-12-1204.

[10] Benjamin MacLellan, Piotr Roztocki, Stefanie Czischek, and Roger G Melko. End-to-end variational quantum sensing. *npj Quantum Information*, 2024. DOI:10.1038/s41534-024-00914-w.

[11] Corentin Lanore, Federico Grasselli, Xavier Valcarce, Jean-Daniel Bancal, and Nicolas Sangouard. Automated generation of photonic circuits for bell tests with homodyne measurements. *arXiv preprint*, 2024. DOI:10.48550/arXiv.2410.19670.

[12] Tommaso Dorigo, Andrea Giammanco, Pietro Vischia, Max Aehle, Mateusz Bawaj, Alexey Boldyrev, Pablo de Castro Manzano, Denis Derkach, Julien Donini, Auralee Edelen, Federica Fanzago, Nicolas R. Gauger, Christian Glaser, Atılım G. Baydin, Lukas Heinrich, et al. Toward the end-to-end optimization of particle physics instruments with differentiable programming. *Reviews in Physics*, 10, 2023. DOI:10.1016/j.revip.2023.100085.

[13] J Rolla, A Machtay, A Patton, W Banzhaf, A Connolly, R Debolt, L Deer, E Fahimi, E Ferstle, P Kuzma, et al. Using evolutionary algorithms to design antennas with greater sensitivity to ultrahigh energy neutrinos. *Physical Review D*, 108(10):102002, 2023. DOI:10.1103/PhysRevD.108.102002.

[14] Giles C Strong, Maxime Lagrange, Aitor Orio, Anna Bordignon, Florian Bury, Tommaso Dorigo, Andrea Giammanco, Mariam Heikal, Jan Kieseler, Max Lamparth, et al. Tomopt: differential optimisation for task- and constraint-aware design of particle detectors in the context of muon tomography. *Machine Learning: Science and Technology*, 5, 2024. DOI:10.1088/2632-2153/ad52e7.

[15] Sean Molesky, Zin Lin, Alexander Y Piggott, Weiliang Jin, Jelena Vucković, and Alejandro W Rodriguez. Inverse design in nanophotonics. *Nature Photonics*, 2018. DOI:10.1038/s41566-018-0246-9.

[16] Neil V Sapra, Ki Youl Yang, Dries Vercruysse, Kenneth J Leedle, Dylan S Black, R Joel England, Logan Su, Rahul Trivedi, Yu Miao, Olav Solgaard, et al. On-chip integrated laser-driven particle accelerator. *Science*, 2020. DOI:10.1126/science.aay5734.

[17] Joshua Yang, Melissa A Guidry, Daniil M Lukin, Kiyoul Yang, and Jelena Vučković. Inverse-designed silicon carbide quantum and nonlinear photonics. *Light: Science & Applications*, 2023. DOI:10.1038/s41377-023-01253-9.

[18] Elias Nehme, Daniel Freedman, Racheli Gordon, Boris Ferdman, Lucien E Weiss, Onit Alalouf, Tal Naor, Reut Orange, Tomer Michaeli, and Yoav Shechtman. Deepstorm3d: dense 3d localization microscopy and psf design by deep learning. *Nature methods*, 17(7), 2020. DOI:10.1038/s41592-020-0853-5.

[19] Carla Rodríguez, Sören Arlt, Leonhard Möckl, and Mario Krenn. Automated discovery of experimental designs in super-resolution microscopy with xlumina. *Nature Communications*, 2024. DOI:10.1038/s41467-024-54696-y.

[20] Mario Krenn, Yehonathan Drori, and Rana X Adhikari. Digital discovery of interferometric gravitational wave detectors. *Phys. Rev. X*, 15, 2025. DOI:10.1103/PhysRevX.15.021012.

[21] Gustavo Mendes Platt, Xin-She Yang, and Antônio José Silva Neto. *Computational intelligence, optimization and inverse problems with applications in engineering.* Springer, 2018. DOI:10.1007/978-3-319-96433-1.

[22] William J Dally, Stephen W Keckler, and David B Kirk. Evolution of the graphics processing unit (gpu). *IEEE Micro*, 41, 2021. DOI:10.1109/MM.2021.3113475.

[23] Cristina Silvano, Daniele Ielmini, Fabrizio Ferrandi, Leandro Fiorin, Serena Curzel, Luca Benini, Francesco Conti, Angelo Garofalo, Cristian Zambelli, Enrico Calore, et al. A survey on deep learning hardware accelerators for heterogeneous hpc platforms. *ACM Computing Surveys*, 2025. DOI:10.1145/3729215.

[24] Takashi Asano and Susumu Noda. Optimization of photonic crystal nanocavities based on deep learning. *Optics express*, 2018. DOI:10.1364/OE.26.032704.

[25] Wei Ma, Zhaocheng Liu, Zhaxylyk A Kudyshev, Alexandra Boltasseva, Wenshan Cai, and Yongmin Liu. Deep learning for the design of photonic structures. *Nature photonics*, 15(2):77–90, 2021. DOI:10.1038/s41566-020-0685-y.

[26] Jiaqi Jiang, Mingkun Chen, and Jonathan A Fan. Deep neural networks

for the evaluation and design of photonic devices. *Nature Reviews Materials*, 6(8): 679–700, 2021. DOI:10.1038/s41578-020-00260-1.

[27] Paweł Kudela, Abdalraheem Ijjeh, Maciej Radzienski, Marco Miniaci, Nicola Pugno, and Wieslaw Ostachowicz. Deep learning aided topology optimization of phononic crystals. *Mechanical Systems and Signal Processing*, 200:110636, 2023. DOI:10.1016/j.ymssp.2023.110636.

[28] Thomas J Grady, Rishi Khan, Mathias Louboutin, Ziyi Yin, Philipp A Witte, Ranveer Chandra, Russell J Hewett, and Felix J Herrmann. Model-parallel fourier neural operators as learned surrogates for large-scale parametric pdes. *Computers & Geosciences*, 178, 2023. DOI:10.1016/j.cageo.2023.105402.

[29] Ziyi Yin, Rafael Orozco, Mathias Louboutin, and Felix J Herrmann. Solving multiphysics-based inverse problems with learned surrogates and constraints. *Advanced Modeling and Simulation in Engineering Sciences*, 10, 2023. DOI:10.1186/s40323-023-00252-0.

[30] Vittorio Peano, Florian Sapper, and Florian Marquardt. Rapid exploration of topological band structures using deep learning. *Physical Review X*, 11(2):021052, 2021. DOI:10.1103/PhysRevX.11.021052.

[31] Kylian Schmidt, Nikhil Kota, Jan Kieseler, Andrea De Vita, Markus Klute, Abhishek, Max Aehle, Muhammad Awais, Alessandro Breccia, Riccardo Carroccio, Long Chen, Tommaso Dorigo, Nicolas R. Gauger, Enrico Lupi, Federico Nardi, Xuan Tung Nguyen, Fredrik Sandin, Joseph Willmore, and Pietro Vischia. End-to-end detector optimization with diffusion models: A case study in sampling calorimeters. *Particles*, 8(2), 2025. DOI:10.3390/particles8020047.

[32] Daniel David Brown, Andreas Freise, Huy Tuong Cao, Alexei Ciobanu, Jeremie Gobeil, Anna Green, Paul Hapke, Philip Jones, Miron van der Kolk, Kevin Kuns, Sean Leavey, Jonathan Warren Perry, Samuel Rowlinson, and Mischa Sallé. Finesse. 2025. DOI:10.5281/zenodo.12662017.

[33] Benjamin P Abbott, Richard Abbott, TD Abbott, MR Abernathy, Fausto Acernese, Kendall Ackley, Carl Adams, Thomas Adams, Paolo Addesso, RX Adhikari, et al. Observation of gravitational waves from a binary black hole merger. *Physical Review Letters*, 116(6), 2016. DOI:10.1103/PhysRevLett.116.061102.

[34] Burr Settles. *Active learning literature survey.* University of Wisconsin-Madison Department of Computer Sciences, 2009.

[35] Dongyuan Li, Zhen Wang, Yankai Chen, Renhe Jiang, Weiping Ding, and Manabu Okumura. A survey on deep active learning: Recent advances and new frontiers. *IEEE Transactions on Neural Networks and Learning Systems*, 2025. DOI:10.1109/TNNLS.2024.3396463.

[36] Kirthevasan Kandasamy, Gautam Dasarathy, Junier B Oliva, Jeff Schneider, and Barnabás Póczos. Gaussian process bandit optimisation with multi-fidelity evaluations. *Advances in neural information processing systems*, 29, 2016. DOI:10.48550/arXiv.1603.06288.

[37] Kevin Tran and Zachary W Ulissi. Active learning across intermetallics to guide discovery of electrocatalysts for co2 reduction and h2 evolution. *Nature Catalysis*, 1, 2018. DOI:10.1038/s41929-018-0142-1.

[38] Turab Lookman, Prasanna V Balachandran, Dezhen Xue, and Ruihao Yuan. Active learning in materials science with emphasis on adaptive sampling using uncertainties for targeted design. *npj Computational Materials*, 2019. DOI:10.1038/s41524-019-0153-8.

[39] Derek van Tilborg and Francesca Grisoni. Traversing chemical space with active deep learning for low-data drug discovery. *Nature Computational Science*, 2024. DOI:10.1038/s43588-024-00697-2.

[40] Alexander Forrester, Andras Sobester, and Andy Keane. *Engineering design via surrogate modelling: a practical guide.* John Wiley & Sons, 2008. DOI:10.1002/9780470770801.

[41] Paolo Testolina, Mattia Lecci, Mattia Rebato, Alberto Testolin, Jonathan Gambini, Roberto Flamini, Christian Mazzucco, and Michele Zorzi. Enabling simulation-based optimization through

machine learning: A case study on antenna design. In *2019 IEEE Global Communications Conference.* IEEE, 2019. DOI:10.1109/globecom38437.2019.9013240.

[42] Kai-Hung Chang and Chin-Yi Cheng. Learning to simulate and design for structural engineering. In *International conference on machine learning,* pages 1426–1436. PMLR, 2020. DOI:10.48550/arXiv.2003.09103.

[43] Sunwoong Yang, Sanga Lee, and Kwanjung Yee. Inverse design optimization framework via a two-step deep learning approach: application to a wind turbine airfoil. *Engineering with Computers,* 39(3):2239–2255, 2023. DOI:10.1007/s00366-022-01617-6.

[44] Wei Ma, Feng Cheng, and Yongmin Liu. Deep-learning-enabled on-demand design of chiral metamaterials. *ACS nano,* 12(6):6326–6334, 2018. DOI:10.1021/acsnano.8b03569.

[45] Qi Wang and Longfei Zhang. Inverse design of glass structure with deep graph neural networks. *Nature communications,* 12 (1):5359, 2021. DOI:10.1038/s41467-021-25490-x.

[46] Dezhen Xue, Prasanna V Balachandran, John Hogden, James Theiler, Deqing Xue, and Turab Lookman. Accelerated search for materials with targeted properties by adaptive design. *Nature communications,* 7(1): 1–9, 2016. DOI:10.1038/ncomms11241.

[47] Cynthia Shen, Mario Krenn, Sagi Eppel, and Alan Aspuru-Guzik. Deep molecular dreaming: inverse machine learning for de-novo molecular design and interpretability with surjective representations. *Machine Learning: Science and Technology,* 2021. DOI:10.1088/2632-2153/ac09d6.

[48] Tareq Jaouni, Sören Arlt, Carlos Ruiz-Gonzalez, Ebrahim Karimi, Xuemei Gu, and Mario Krenn. Deep quantum graph dreaming: deciphering neural network insights into quantum experiments. *Machine Learning: Science and Technology,* 2024. DOI:10.1088/2632-2153/ad2628.

[49] Auralee Edelen, Nicole Neveu, Matthias Frey, Yannick Huber, Christopher Mayes, and Andreas Adelmann. Machine learning for orders of magnitude speedup in multiobjective optimiza-

tion of particle accelerator systems. *Phys. Rev. Accel. Beams,* 23:044601, 2020. DOI:10.1103/PhysRevAccelBeams.23.044601.

[50] Sergey Shirobokov, Vladislav Belavin, Michael Kagan, Andrei Ustyuzhanin, and Atilim Gunes Baydin. Black-box optimization with local generative surrogates. *Advances in neural information processing systems,* 33, 2020. DOI:10.48550/arXiv.2002.04632.

[51] Kinga Anna Wozniak, Stephen Mulligan, Jan Kieseler, Markus Klute, Francois Fleuret, and Tobias Golling. End-to-end optimal detector design with mutual information surrogates. *arXiv preprint,* 2025. DOI:10.48550/arXiv.2503.14342.

[52] Dianjing Liu, Yixuan Tan, Erfan Khoram, and Zongfu Yu. Training deep neural networks for the inverse design of nanophotonic structures. *ACS Photonics,* 5(4):1365–1369, 2018. DOI:10.1021/acsphotonics.7b01377.

[53] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *Advances in Neural Information Processing Systems,* 2020. DOI:10.48550/arXiv.2006.10739.

[54] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics,* 378:686–707, 2019. DOI:10.1016/j.jcp.2018.10.045.

[55] George Em Karniadakis, Ioannis G. Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics,* 3(6), 2021. DOI:10.1038/s42254-021-00314-5.

[56] Lu Lu, Pengzhan Jin, Guofei Pang, Zhongqiang Zhang, and George Em Karniadakis. Learning nonlinear operators via deeponet based on the universal approximation theorem of opera-

tors. *Nature Machine Intelligence*, 3, 2021. DOI:10.1038/s42256-021-00302-5.

[57] Kamyar Azizzadenesheli, Nikola Kovachki, Zongyi Li, Miguel Liu-Schiaffini, Jean Kossaifi, and Anima Anandkumar. Neural operators for accelerating scientific simulations and design. *Nature Reviews Physics*, 2024. DOI:10.1038/s42254-024-00712-5.

[58] Mohammadamin Mahmoudabadbozchelou, George Em Karniadakis, and Safa Jamali. nn-pinns: Non-newtonian physics-informed neural networks for complex fluid modeling. *Soft Matter*, 18, 2022. DOI:10.1039/D1SM01298C.

[59] Benedikt Alkin, Andreas Fürst, Simon Schmid, Lukas Gruber, Markus Holzleitner, and Johannes Brandstetter. Universal physics transformers: A framework for efficiently scaling neural operators. *Advances in Neural Information Processing Systems*, 2024. DOI:10.48550/arXiv.2402.12365.

[60] Vivek Oommen, Khemraj Shukla, Saaketh Desai, Rémi Dingreville, and George Em Karniadakis. Rethinking materials simulations: Blending direct numerical simulations with neural operators. *npj Computational Materials*, 10, 2024. DOI:10.1038/s41524-024-01319-1.

[61] Sandeep Suresh Cranganore, Andrei Bodnar, Arturs Berzins, and Johannes Brandstetter. Einstein fields: A neural perspective to computational general relativity. *arXiv preprint*, 2025. DOI:10.48550/arXiv.2507.11589.

[62] Sanjoy Dasgupta. Analysis of a greedy active learning strategy. In *Proceedings of the 18th International Conference on Neural Information Processing Systems*, NIPS'04, Cambridge, MA, USA, 2004. MIT Press. DOI:10.5555/2976040.2976083.

[63] Steve Hanneke. A bound on the label complexity of agnostic active learning. In *Proceedings of the 24th International Conference on Machine Learning*, ICML '07, New York, NY, USA, 2007. Association for Computing Machinery. DOI:10.1145/1273496.1273541.

[64] Ravi Ganti and Alexander Gray. Upal: Unbiased pool based active learning. In *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research. PMLR, 2012. DOI:10.48550/arXiv.1111.1784.

[65] Alon Gonen, Sivan Sabato, and Shai Shalev-Shwartz. Efficient active learning of halfspaces: an aggressive approach. *J. Mach. Learn. Res.*, 2013. DOI:10.5555/2567709.2567744.

[66] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011. DOI:10.48550/arXiv.1011.0686.

[67] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 2017. DOI:10.48550/arXiv.1612.01474.

[68] Keze Wang, Dongyu Zhang, Ya Li, Ruimao Zhang, and Liang Lin. Cost-effective active learning for deep image classification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2017. DOI:10.1109/TCSVT.2016.2589879.

[69] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *International Conference on Learning Representations*, 2018. DOI:10.48550/arXiv.1708.00489.

[70] Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling estimation of error reduction. In *In Proc. 18th International Conf. on Machine Learning*, 2001.

[71] Ashish Kapoor, Kristen Grauman, Raquel Urtasun, and Trevor Darrell. Active learning with gaussian processes for object categorization. In *2007 IEEE 11th International Conference on Computer Vision*, 2007. DOI:10.1109/ICCV.2007.4408844.

[72] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: representing model uncertainty in deep learning. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, ICML'16, 2016. DOI:10.48550/arXiv.1506.02142.

[73] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning*, ICML'17, 2017. DOI:10.48550/arXiv.1703.02910.

[74] Raphaël Pestourie, Youssef Mroueh, Thanh V. Nguyen, Payel Das, and Steven G. Johnson. Active learning of deep surrogates for pdes: application to metasurface design. *npj Computational Materials*, 6, 2020. DOI:10.1038/s41524-020-00431-2.

[75] Sudhanshu Singh, Rahul Kumar, Praveen Singh, and Ravi Hegde. Active learning for efficient nanophotonics inverse design in large and diverse design spaces. *Opt. Express*, 2025. DOI:10.1364/OE.559669.

[76] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint*, 2014. DOI:10.48550/arXiv.1412.6572.

[77] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representations*, 2017. DOI:10.48550/arXiv.1611.01236.

[78] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. DOI:10.48550/arXiv.1706.06083.

[79] Albert Einstein. Näherungsweise integration der feldgleichungen der gravitation. *Sitzungsberichte der Königlich Preußischen Akademie der Wissenschaften*, pages 688–696, 1916.

[80] Sukin Sim, Peter D Johnson, and Alán Aspuru-Guzik. Expressibility and entangling capability of parameterized quantum circuits for hybrid quantum-classical algorithms. *Advanced Quantum Technologies*, 2019. DOI:10.1002/qute.201900070.

[81] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*, 2014. DOI:10.48550/arXiv.1412.6980.

[82] Geoffrey E Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 2002.

[83] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.

[84] Caglar Gulcehre, Marcin Moczulski, Misha Denil, and Yoshua Bengio. Noisy activation functions. In *International conference on machine learning*. PMLR, 2016. DOI:10.48550/arXiv.1603.00391.

[85] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. DOI:10.48550/arXiv.2010.11929.

[86] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34, 2021. DOI:10.48550/arXiv.2105.01601.

[87] Xiaoxin He, Bryan Hooi, Thomas Laurent, Adam Perold, Yann Lecun, and Xavier Bresson. A generalization of ViT/MLP-mixer to graphs. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*. PMLR, 2023. DOI:10.48550/arXiv.2212.13350.

[88] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *Proceedings of the 36th International Conference on Machine Learning*. PMLR, 2019. DOI:10.48550/arXiv.1806.08734.

[89] Ronen Basri, Meirav Galun, Amnon Geifman, David Jacobs, Yoni Kasten, and Shira Kritchman. Frequency bias in neural networks for input of non-uniform density. In *Proceedings of the 37th International Con-*

*ference on Machine Learning*. PMLR, 2020. DOI:10.48550/arXiv.2003.04560.

[90] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, 2017. DOI:10.48550/arXiv.1706.03762.

[91] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *Proceedings of the 37th International Conference on Machine Learning*. PMLR, 2020. DOI:10.48550/arXiv.2002.04745.

[92] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*. PMLR, 2019. DOI:10.48550/arXiv.1810.00825.

[93] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019. DOI:10.18653/v1/N19-1423.

[94] Philippe Schwaller, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A Hunter, Costas Bekas, and Alpha A Lee. Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS central science*, 2019. DOI:10.1021/acscentsci.9b00576.

[95] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 2021. DOI:10.1038/s41586-021-03819-2.

[96] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021. DOI:10.1609/aaai.v35i12.17325.

[97] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 2021. DOI:10.48550/arXiv.2106.01345.

[98] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. *Advances in neural information processing systems*, 2021. DOI:10.48550/arXiv.2106.02039.

[99] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. DOI:10.48550/arXiv.2103.14030.

[100] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on pattern analysis and machine intelligence*, 2022. DOI:10.1109/TPAMI.2022.3152247.

[101] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML '09. Association for Computing Machinery, 2009. DOI:10.1145/1553374.1553380.

## A UIFO characterization

### A.1 Consistent distances

In the UIFO framework, the optical components are located within an irregular rectangular grid (see Figure 5). To locate each element, we set the distances between columns and rows of beam splitters and then locate the mirrors in between based on fractions of the largest distances. We fix the distances between the outer mirrors and their closest light source (or detector), since they do not affect the GWD performance. The values of the vertical and horizontal distances, $V_i$ and $H_j$, are also fractions of a larger tunable length, which can take up to 4km.

With the fraction based location the mirrors could overlap. To prevent that nonphysical situation, we impose a minimum value $\epsilon$ to all the distances described. This works because the number of elements between beam splitters is always the same, and therefore the number of distances. If we wanted to remove elements we should be cautious with these fixed gaps.
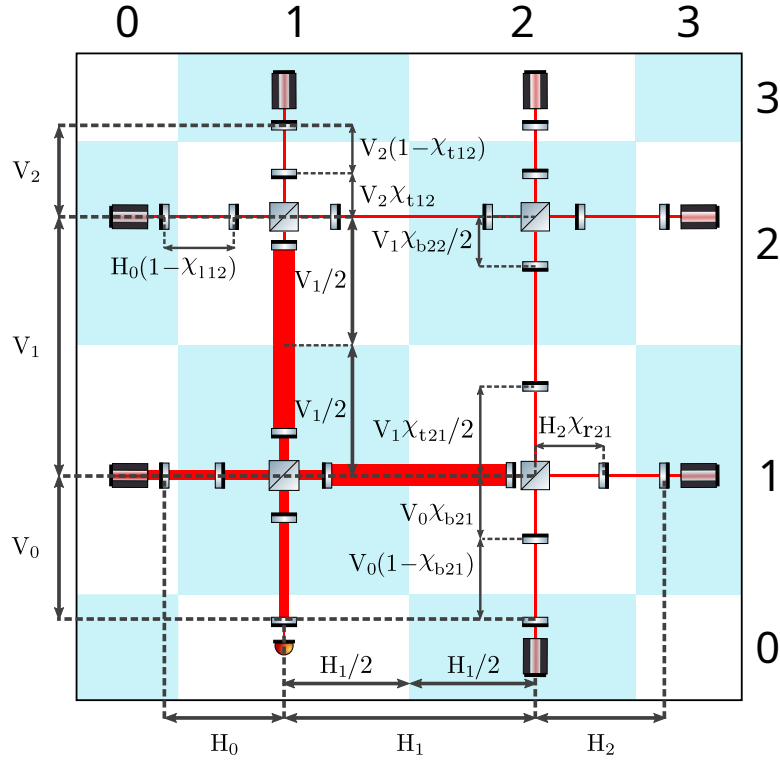


Figure 5: **Positioning of the optical elements.** To keep the grid structure, we locate the mirrors of the optical path based on fractions of larger distances, vertical $\{V_i\}$ and horizontal $\{H_j\}$, which are also fractions of a maximum vertical and horizontal size, up to 4km. For the sake of readability, we keep the mirrors within the cells, so the mirrors between beam splitters can be localized as a fraction of half the gap between beam splitters. The distance between light sources and their nearest mirrors is fixed, since that does not affect the design performance. Same applies for the detector and its nearest mirror.

### A.2 Discrete design choices

Having set the UIFO topology and distances there are a few discrete design choices to make, like where to locate the detector, or which type of light sources to use. We also have to chose if we place beam splitters or Faraday isolators in the grid intersections, and how to orient them. We could also decide, for each optical component, whether we want it to be there at all.

In our experiments, we always place the detector in the same grid cell and only use lasers as light sources, not squeezers. We do not use Faraday isolators either and the orientation of the beam splitters is always the same. We are keeping it simple for now.

## A.3  Optical components

Except for the detector, all optical components have, at least, one continuous parameter to tune. Table 3 summarizes the range of parameter values for the UIFO designs, before they are normalized (from 0 to 1) for training. Finesse can also simulate other elements and parameters, all listed and documented [32].

It is worth noticing that when you 'turn off' a light source, both for squeezers and lasers, by setting their parameters to 0, their effect on the system ceases. That is not the case for mirrors and beam splitters. We can never ignore them because they have a minimum loss, and because they are still susceptible to power damages. To properly 'turn them off' we must remove them completely.

| Components | Parameters | Min | Max |
|---|---|---|---|
| Mirrors and beam splitters | Reflectivity | 0 | 0.99998 |
| | Phase (°) | 0 | 360 |
| Lasers | Power (W) | 0 | 200 |
| Squeezers | Factor (db) | 0 | 20 |
| | Angle (°) | 0 | 360 |

Table 3: Range of values (before normalization) for different optical components. In Finesse, the mass of mirrors and beam splitters is infinite by default, we leave it like that in all our experiments. Mirrors and beam splitters have a minimum loss and transmission of 5ppm and 15ppm, respectively, which limits their reflectivity.

## B  Datasets

The 1 million training samples and 100K testing samples are stored in H5 data files with 1000 pairs of inputs-outputs each, normalized from 0 to 1. The inputs are the linearly normalized parameters of distances (lengths up to 4km and fractions of these) and optical components (described in Table 3). The loss and masses of mirrors and beamsplitters are not included as they are constant.

The outputs are the sensitivity values of 101 frequency values (a geometric sequence from 1Hz to 10kHz), and the light power acting on the detector, mirrors, and beam splitters. As with the loss and masses, we do not include the frequency values in the dataset, as they are always the same. Since both sensitivities and powers range across several orders of magnitude, we must not rescale them directly from 0 to 1. Instead, and only during training/testing, we rescale the logarithm of their values

$$\mathcal{N}(x) \equiv \frac{\log_{10}(x + \epsilon) + b}{S}, \tag{3}$$

where the values of $e$, $b$, and $S$ are (0, 30, 30) for the sensitivities, and $(10^{-6}, 6, 16)$ for the powers. These choices are motivated by the expected range of sensitivities, usually from $10^{-25}$ to $10^{-10}$ (by definition, never 0 or higher than 1), and the expected range of powers, from 0 to $10^7$ (beyond that, everything burns).

## C  Surrogate model architecture

To mimic the behavior of the Finesse simulator we trained a transformer-based surrogate model grids [86, 87]. The parameters of neighbor elements were merged into tokens. We divided the grid by the beamsplitters (see Figure 6) and add their parameters to the four neighboring subgraphs. These overlapping patches were projected into a higher-dimensional space via a linear layer followed by a sinusoidal activation. This transformation, also known as Fourier feature mapping [53], proved essential for capturing the high-frequency features of the GWD (see Figure 1), overcoming the spectral

bias of standard neural networks [88, 89]. The tokens are then processed by 8 transformer layers. After processing the tokens, they were projected into a lower dimension and concatenated. With a final multilayer perceptron, we obtain an array of predicted sensitivities and optical powers.

Since the transformer encoder is invariant under permutation of patches, we add 2 trainable positional encoding to the patches: before and after they are processed by the Fourier mapping.
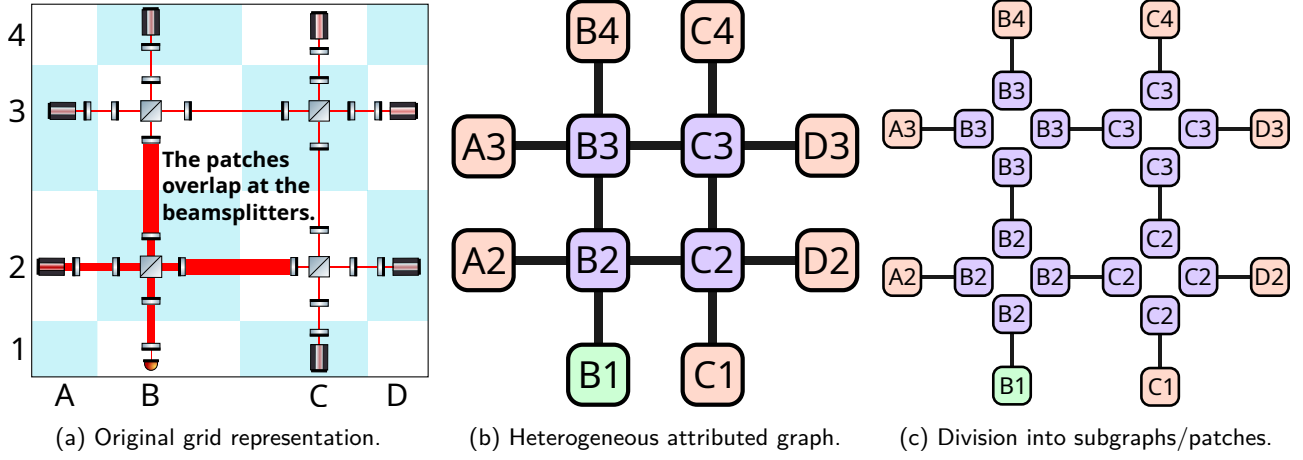


(a) Original grid representation.   (b) Heterogeneous attributed graph.   (c) Division into subgraphs/patches.

Figure 6: **Patch extraction from the gravitational wave detector.** We encoded GWD into rectangular grids **(a)**, which can be seen as graphs with different types of nodes, and multiple attributes at nodes and edges **(b)**. Dividing the rectangular graph into overlapping subgraphs, we group the parameters that characterize them into patches **(c)**. The three types of nodes, light sources (orange), beam splitters (purple), and detector (green), while the mirror properties are encoded into the edges that connect the nodes. Since the patches have different number of attributes we pad the patches up to a common length. For simplicity, we group the vertical and horizontal distances, $\{V_i\}$ and $\{H_j\}$, into 2 additional patches.

# D   GPU resources used for surrogate-based inverse design

The training times showed in Table 2 are averaged from the data shown in Table 4. The inverse design of 20K samples over 16384 optimization steps took 16, 40 and 85 minutes on average, for the 3 grid sizes. We used *AMD Instinct MI300A APU* units.

| Generation of surrogate | Seed | UIFO size (# parameters) | | |
|:---:|:---:|:---:|:---:|:---:|
| | | 1x1 (29) | 2x2 (85) | 3x3 (169) |
| 0th | 0 | 20' 12" | 27' 00" | 47' 00" |
| | 1 | 20' 05" | 31' 23" | 40' 57" |
| 1st | 0 | 45' 09" | 56' 04" | 1h 30' 32" |
| | 1 | 44' 29" | 56' 44" | 1h 13' 20" |
| 2nd | 0 | 1h 07' 12" | 1h 36' 27" | 2h 25' 46" |
| | 1 | 1h 15' 58" | 1h 47' 48" | 1h 39' 36" |
| 3rd | 0 | 1h 42' 31" | 2h 04' 59" | 3h 06' 26" |
| | 1 | 1h 43' 32" | 2h 25' 05" | 3h 12' 27" |
| 4th | 0 | 2h 20' 13" | 2h 34' 05" | 4h 09' 31" |
| | 1 | 2h 12' 22" | 2h 48' 26" | 3h 25' 32" |
| **Total time** | 0 | 6h 15' 17" | 7h 38' 35" | 11h 59' 15" |
| | 1 | 6h 16' 26" | 8h 29' 26" | 10h 11' 52" |

Table 4: **Training times of the surrogate models for the 3 UIFO sizes.**

# E  Optimization with Finesse

Figure 7 shows the detailed results of the design optimization run with Finesse for 5 days in 256 parallel jobs for each size of the UIFO.
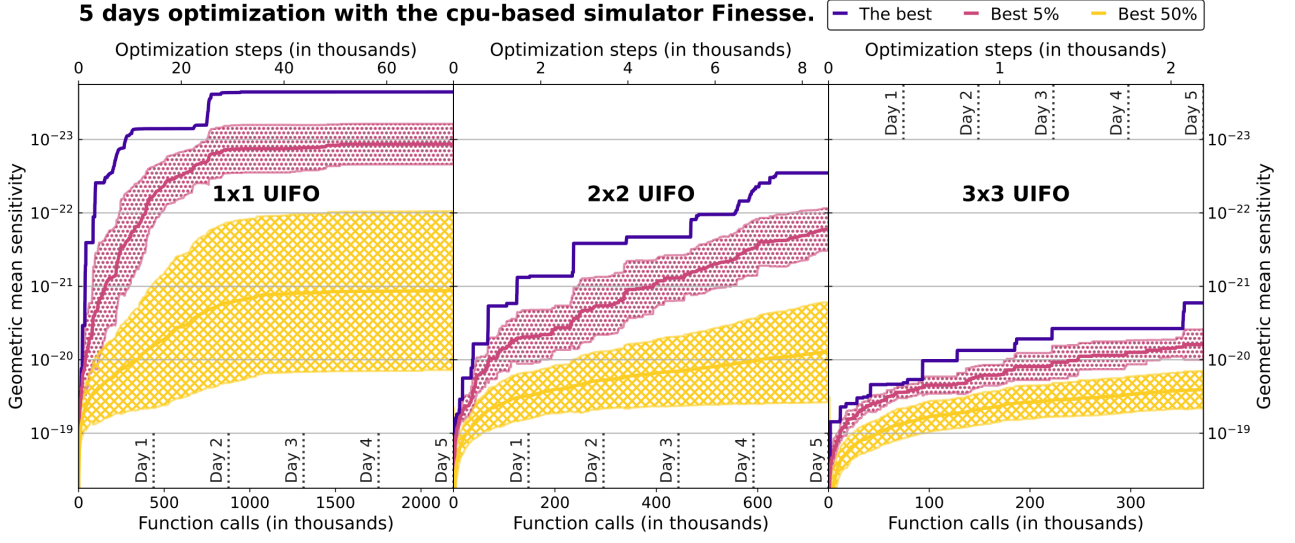


Figure 7: **Optimization of 256 designs with Finesse during 5 days.** Using the Adam optimizer with annealed noise (same than for the surrogates) we optimize 256 independent input parameter sets, randomly initialized, to minimize the loss from Eq.(2). Lacking automatic differentiation, we must use finite differences to compute the gradients, taking 30, 86, and 170 evaluations per gradient for the different UIFO sizes: 1x1, 2x2, and 3x3. With approximate running times of 0.2, 0.6, and 1.2 seconds respectively, each optimization step might take more than 3 minutes for the largest system. The plotted data considers the best sensitivity ever achieved by the samples at a given optimization step, ignoring the decreases from the gradient computations and/or the optimization process.