

# AI Assisted Detector Design for EIC

Karthik Suresh, for AID(2)E Collaboration  
Department of Data Science  
College of William and Mary



WILLIAM & MARY

CHARTERED 1693



# Outline

- Multi Objective Optimization
- Need for AI in detector design – The AID(2)E project
- Closure Test and Project workflows
- Selected works and future studies

# Multi Objective Optimization

Design space spanned by 'x'

$$\min / \max f_m(\mathbf{x}), m = 1, \dots, M$$

$$\text{s.t. } g_j(\mathbf{x}) \leq 0, j = 1, \dots, J$$

$$h_k(\mathbf{x}) = 0, k = 1, \dots, K$$

$$x_i^L \leq x_i \leq x_i^U, i = 1, \dots, N$$

Guo, Kai, et al. *Materials* 8.4 (2021): 1153-1172.

ML method	Characterization	Example applications in mechanical materials design
Linear regression	Model the linear or polynomial relationship between input and output variables	Modulus <sup>101</sup> or strength <sup>102</sup> prediction
Polynomial regression		
Support vector machine; SVM	Separate high-dimensional data space with one or a set of hyperplanes	Strength <sup>103</sup> or hardness <sup>104</sup> prediction; structural topology optimization <sup>105</sup>
Random forest	Construct multiple decision trees for classification or prediction	Modulus <sup>106</sup> or toughness <sup>107</sup> prediction
Feedforward neural network (FNN); MLP	Connect nodes (neurons) with information flowing in one direction	Prediction of modulus, <sup>103,108</sup> strength, <sup>109</sup> toughness <sup>106</sup> or hardness <sup>110</sup> ; prediction of hyperelastic or plastic behaviors, <sup>111-114</sup> identification of collision hot conditions, <sup>115</sup> design of optimal intermetallic <sup>116</sup>
CNNs	Capture features at different hierarchical levels by calculating convolutions, operate on pixel-based or voxel-based data	Prediction of strain fields <sup>103,109</sup> or elastic properties <sup>103,109</sup> of high-concentr composite, modulus of undirectional composites, <sup>117</sup> stress fields in condensed structures, <sup>118</sup> or yield strength of additive-manufactured metals, <sup>119</sup> prediction of fatigue crack propagation in polycarbonate alloys, <sup>120</sup> prediction of crystal plasticity, <sup>121</sup> design of structure composites, <sup>122</sup> design of stretchable graphene <sup>123,124</sup>
Recurrent neural network (RNN); LSTM; GRU	Connect nodes (neurons) forming a directed graph with history information stored in hidden states, operate on sequential data	Prediction of fracture patterns in crystalline solids, <sup>125</sup> prediction of plastic behaviors in heterogeneous materials, <sup>126,127</sup> multi-scale modeling of porous media <sup>128</sup>
Generative adversarial networks (GANs)	Train two opponent neural networks to generate and discriminate approaches, used the new generated data to optimize, generate new data according to the distribution of training set	Prediction of modulus distribution by solving inverse elasticity problems, <sup>129</sup> prediction of stress or stress fields in composites, <sup>130</sup> supersonic design, <sup>131</sup> structural topology optimization, <sup>132</sup> substructured materials design <sup>133</sup>
Gaussian process regression (GPR); Bayesian learning	Treat parameters as random variables and calculate the probability distribution of these variables, quantify the uncertainty of model prediction	Modulus <sup>134</sup> or strength <sup>135</sup> prediction; design of supercompressible and assemble metamaterials <sup>136</sup>
Active learning	Interacts with a user on the fly for labeling new data, augment training data with post-hoc experiments or simulations	Strength prediction <sup>137</sup>
Genetic or evolutionary algorithms	Mimic evolutionary rules for optimizing objective function	Hardware prediction, <sup>138</sup> design of active materials, <sup>139</sup> design of modular metamaterials <sup>140</sup>
Reinforcement learning	Maximize cumulative awards with agents reacting to the environment	Deriving microstructure-based traction- $\sigma$ - $\epsilon$ laws <sup>141</sup>
Graph neural networks (GNNs)	Operate on non-Euclidean data structures, applicable tasks include link prediction, node classification and graph classification	Hardware prediction, <sup>137</sup> architected materials design <sup>142</sup>

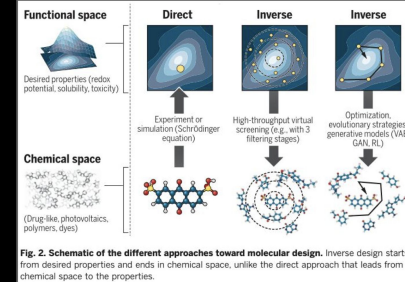


Fig. 2. Schematic of the different approaches toward molecular design. Inverse design starts from desired properties and ends in chemical space, unlike the direct approach that leads from chemical space to the properties.

B. Sanchez-Lengeling, A. Aspuru-Guzik. *Science* 361.6400 (2018): 360-365.

Multiobjective genetic algorithm approach to optimize beam matching and beam transport in high-intensity hadron linacs

M. Yarmohammadi Satri,<sup>1,2,\*</sup> A. M. Lombardi,<sup>3</sup> and F. Zimmermann<sup>3</sup>

<sup>1</sup>School of Particles and Accelerators, Institute for Research in Fundamental Sciences (IPM), P.O. Box 19395-5531, Tehran, Iran  
<sup>2</sup>CERN, 1211 Geneva 23, Switzerland



# Multi Objective Optimization

Design space spanned by 'x'



$$\min / \max \mathbf{f}_m(\mathbf{x}), m = 1, \dots, M$$

Objectives to optimize

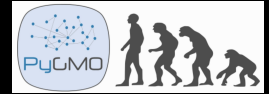
$$\text{s.t. } \mathbf{g}_j(\mathbf{x}) \leq 0, j = 1, \dots, J$$

$$\mathbf{h}_k(\mathbf{x}) = 0, k = 1, \dots, K$$

$$x_i^L \leq x_i \leq x_i^U, i = 1, \dots, N$$

# Multi Objective Optimization

Design space spanned by 'x'



$$\min / \max \mathbf{f}_m(\mathbf{x}), m = 1, \dots, M$$

$$\text{s.t. } \mathbf{g}_j(\mathbf{x}) \leq 0, j = 1, \dots, J$$

Constraints

$$\mathbf{h}_k(\mathbf{x}) = 0, k = 1, \dots, K$$

$$x_i^L \leq x_i \leq x_i^U, i = 1, \dots, N$$

# Multi Objective Optimization

Design space spanned by 'x'

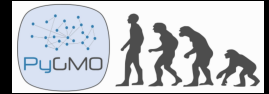
$$\min / \max \mathbf{f}_m(\mathbf{x}), m = 1, \dots, M$$

$$\text{s.t. } \mathbf{g}_j(\mathbf{x}) \leq 0, j = 1, \dots, J$$

$$\mathbf{h}_k(\mathbf{x}) = 0, k = 1, \dots, K$$

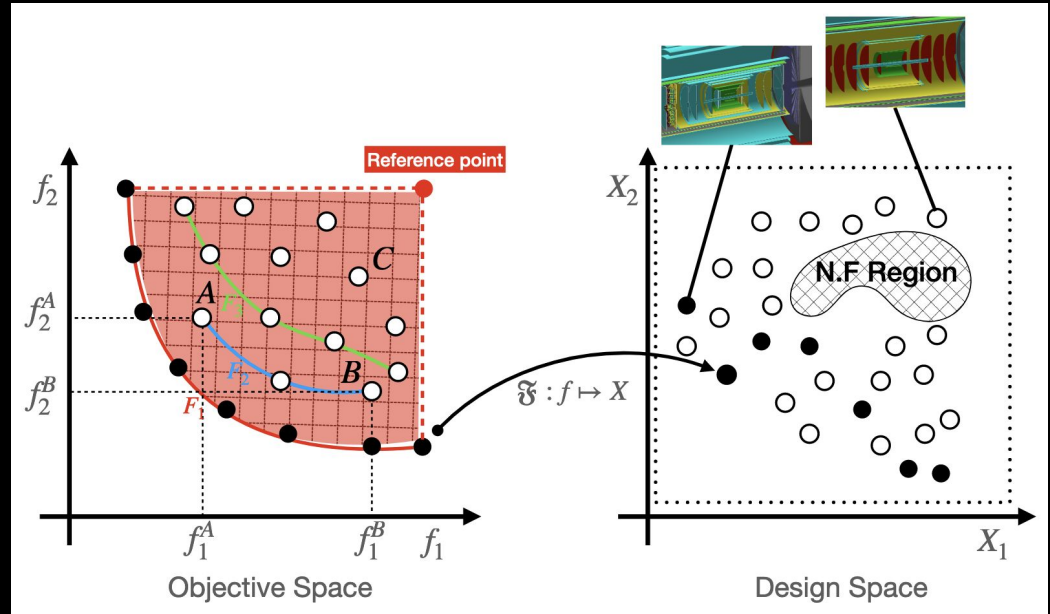
$$x_i^L \leq x_i \leq x_i^U, i = 1, \dots, N$$

Bounded Design Space



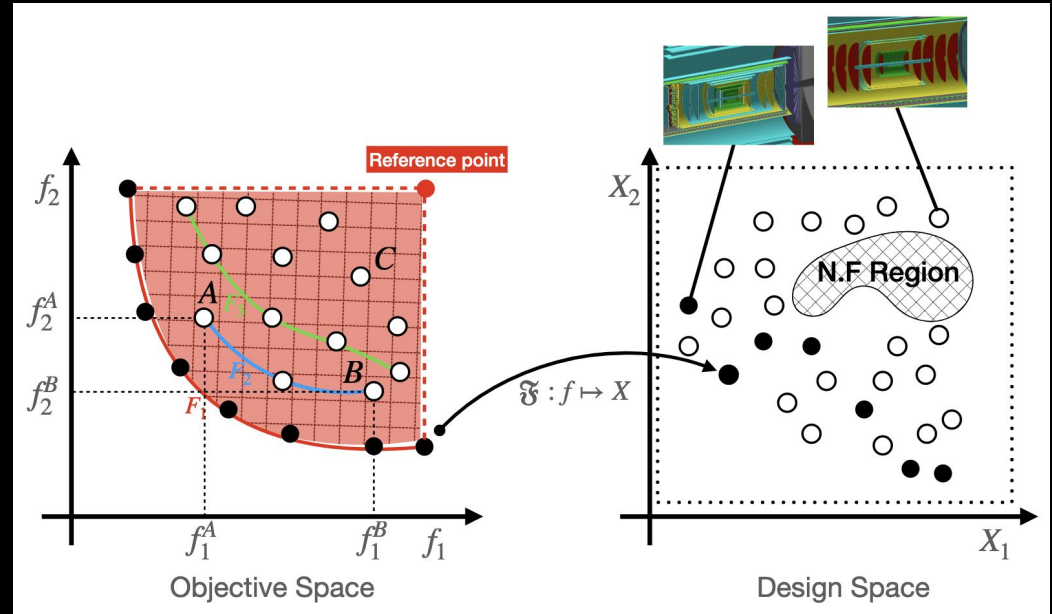
# Multi Objective Optimization : Visual Intro

- Multiple “objectives”
  - Momentum resolution
  - $\theta$  resolution
  - KF efficiency
  - projected  $\theta$  resolution @ PID
- Goal : “Optimize” these Objectives
- Map: “Design” space – “Objective” Space
- Non-Feasible region to be avoided



# Multi Objective Optimization : Visual Intro

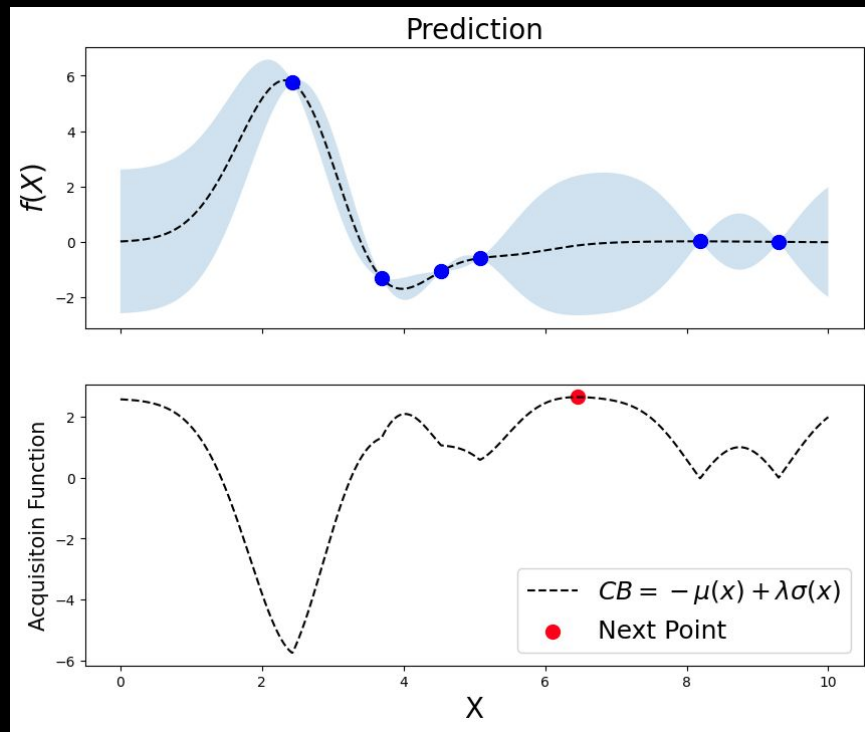
- What is “Optimal”?
  - Non-dominated (Pareto) Solutions
- How to rank solutions?
  - “Fronts” of solutions
- Methods of MOO
  - Evolutionary
  - Bayesian
  - Preferential Learning, etc.





# Multi Objective Optimization through surrogate modelling

- Surrogate Model – A model that will be able to successfully approximate the true function.
- Acquisition Model – A quick evaluator to choose the next point to be computed
  - Based on Exploration and Exploitation in the search space.
  - Critically important, since, this is key in convergence.



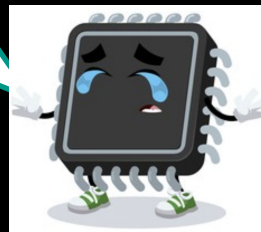
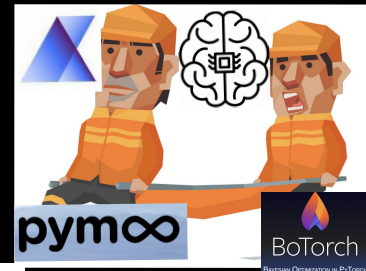
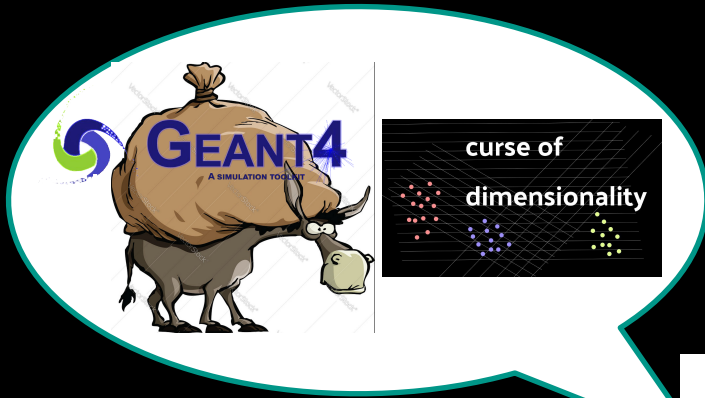
# Large Scale Experiments : An Ideal MOO problem

GEANT4 – computation intensive.

Curse of dimensionality due to multiple Objectives and multidimensional design space

Each Design point requires multiple physics studies and hence increased computational needs

Estimated simulation requirements based on observed performance in 2021.  
<https://arxiv.org/pdf/2205.08607.pdf>



Year	Number of Events [ $\times 10^6$ ]	Storage [TB]	CPU-core hours [Mcore-hrs]
2022	200	50	45
2023 - 2024	100	25	22.5
2025 - 2028	50	12.5	11
2029 - 2030	500	125	110
<b>Total</b>	<b>1600</b>	<b>400</b>	<b>354</b>

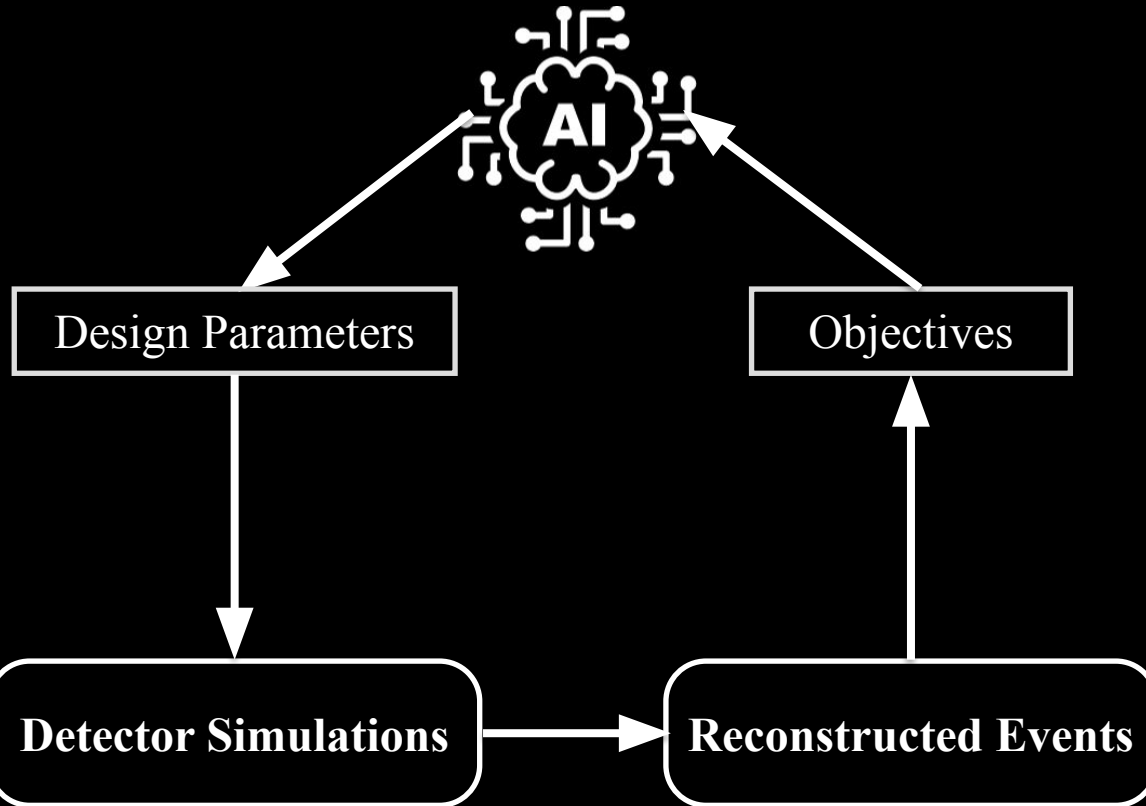
# Workflow for AI Assisted detector design

Benefits from rapid turnaround time from simulations to analysis of high-level reconstructed observables

The EIC SW stack offers multiple features that facilitate AI-assisted design (e.g., modularity of simulation, reconstruction, analysis, easy access to design parameters, automated checks, etc.)

Leverages heterogeneous computing

Need to develop end to end pipeline

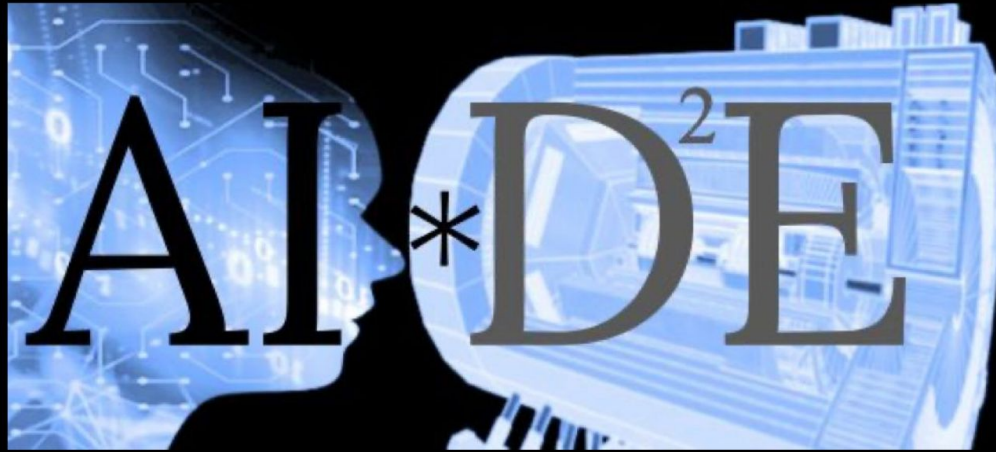


Desired kinematic range

## AID(2)E: AI-Assisted Detector Design at EIC



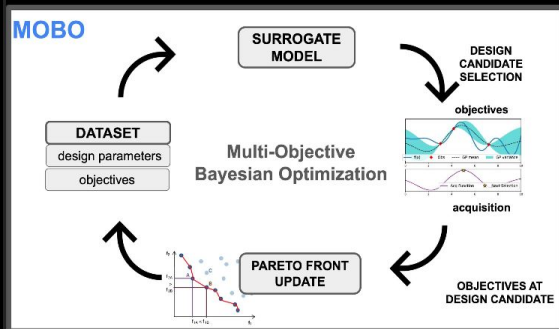
DE-FOA-0002785



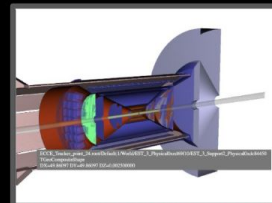
BNL, T. Wenaus  
CUA, T. Horn  
Duke, A. Vossen  
JLab, M. Diefenthaler  
W&M, CF



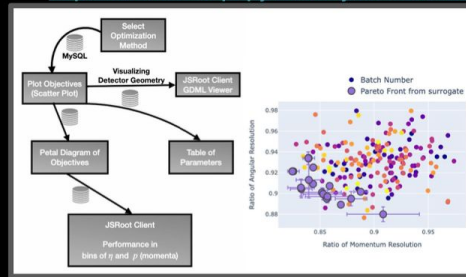
# The AID(2)E Project



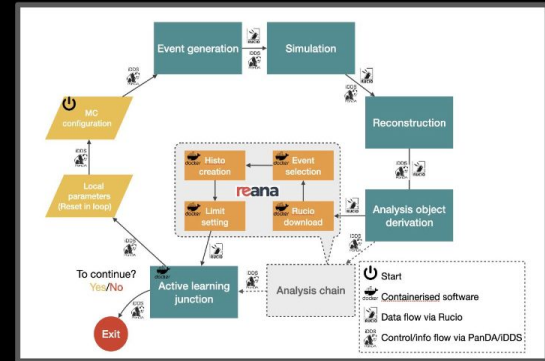
(i) Will contribute to advance the boundaries of MOBO complexity to accommodate a large number of objectives and will explore usage of physics-inspired approaches



<https://ai4eicdetopt.pythonanywhere.com/>



(ii) Development of suite of data science tools for interactive navigation of Pareto front (multi-dim design with multiple objectives)



(iii) Will leverage cutting-edge workload management systems capable of operating at massive data and handle complex workflows

Examining solutions on the Pareto front of ePIC at different values of the budget can have great cost benefits

A fractional improvement in the objectives translates to a more efficient use of beam time which will make up a majority of the cost of the EIC over its lifetime

# Project Workflow

AID(2)E Wrapper

## AID(2)E Pipeline Thrusts of development

Distributed  
Computing

Optimizer

eg. MOBO Algorithm

ePIC Software  
Heavy simulations

Computations  
during an iteration

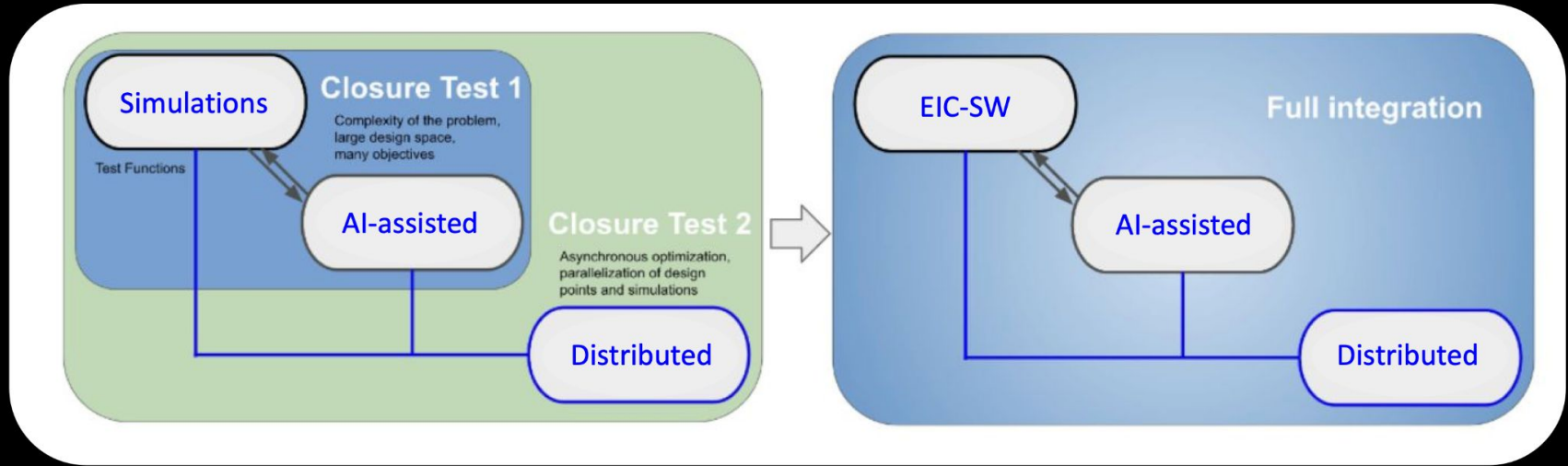
Total number of  
iterations to converge

Detector  
Simulations  
EIC-SIM-RECON

EIC Analysis  
Physics/Detector  
response



# Project workflow



# Closure Test 1 – Stress testing SoTA MOBO

## Gaussian Process $O(n^3)$

- The PDF prior distribution, that describes the Design space to objective. This is the surrogate model.
- SAAS<sup>[1]</sup> priors have been proven to be successful upto 388 design dimensions.
- Assumes several design variables has increased importance compared to others
- Computational expensive as iteration increases
- Benefit from GPU hardware acceleration

## Bayesian Sampling from posteriors NUTS – $O(Md^{5/4})$ <sup>[NUTS]</sup>

- Sample L points from the posterior distribution.
- HMC is a popular algorithm
- Mainly depends on the Number of objectives and design space dimensions
- Has minimal dependence on iteration.
- GPU acceleration through JAX backend.

## Acquisition function qNEHVI – $O(M(n+i)^M)$ <sup>[2]</sup>

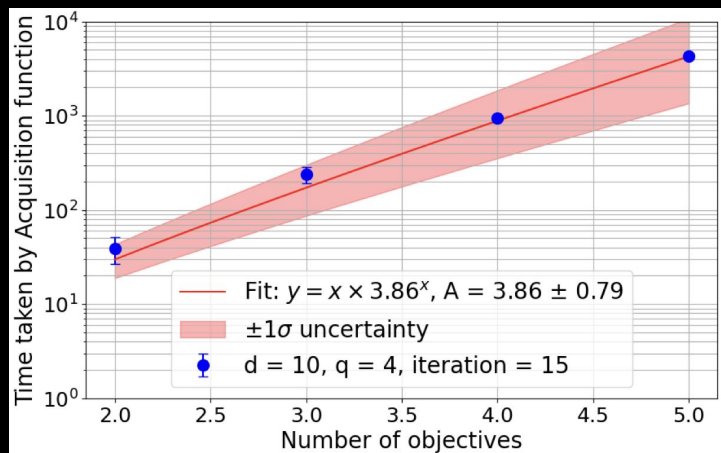
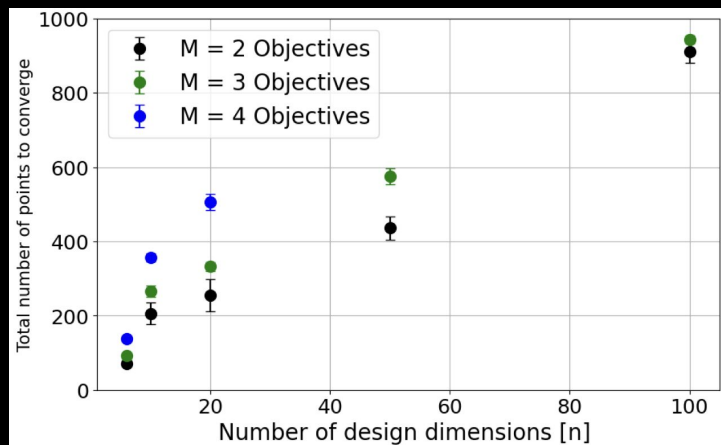
- A “cheaper” function to evaluate as a proxy for the black box function
- Identifies points of maximum improvements hence, the name
- Scales nonlinearly with iteration, total points explored, design space and objective space.
- Partially benefitted by GPU acceleration.



# Closure Test 1 – Stress testing MOBO

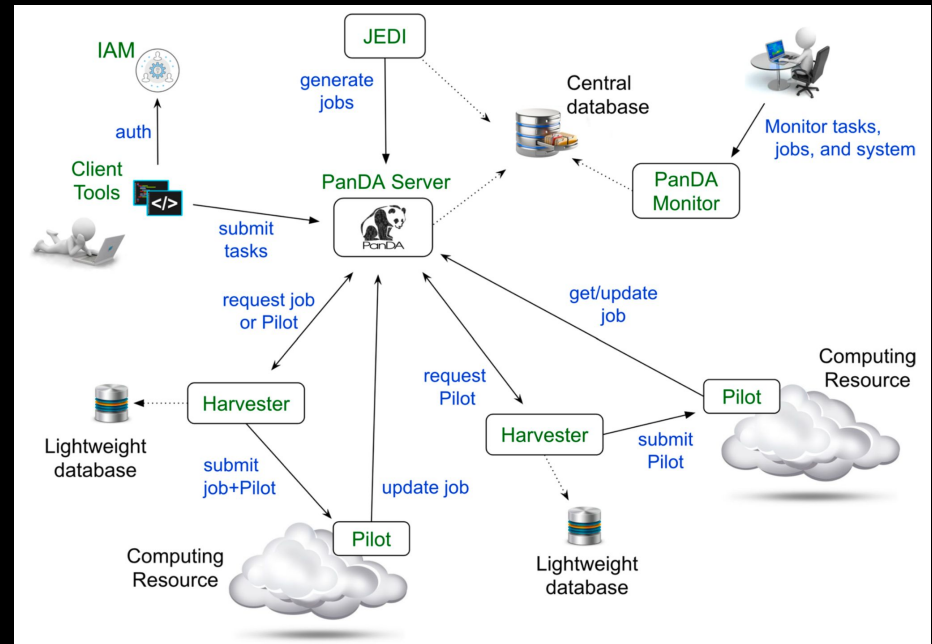
- Stress test the SoTA algorithm used for optimization
- MOBO stress-testing for problems with increasing complexity (design and objectives) and known Pareto

[arXiv:2405.16279](https://arxiv.org/abs/2405.16279)



# Closure Test 2: PanDA/iDDS integration

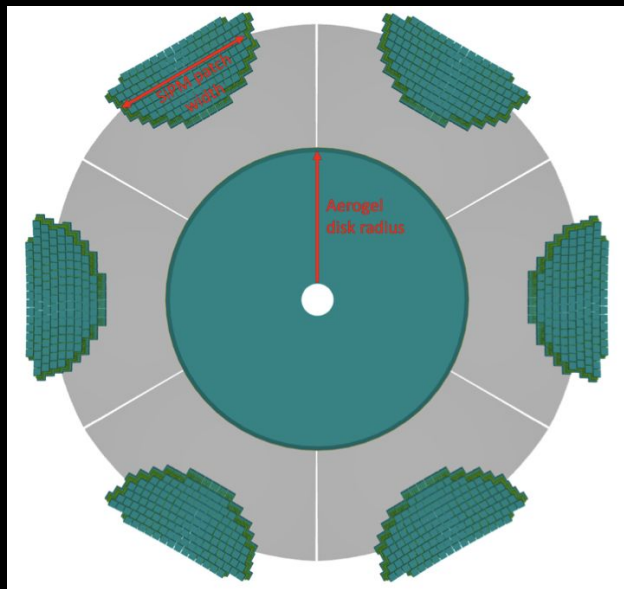
- Stress test scalability across distributed resources
- Integrate PanDA/iDDS AI/ML service to support MOBO workflow for design optimization



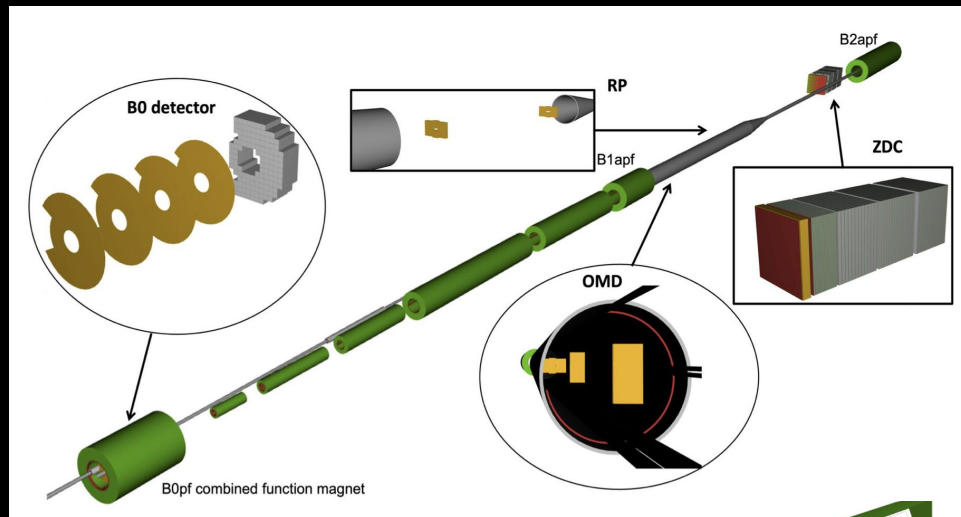
[PanDA: Production and Distributed Analysis System. \*Comput Softw Big Sci\* 8, 4 \(2024\)](#)

# Current Detector Subsystems for optimization in ePIC

## d-RICH detector at EIC

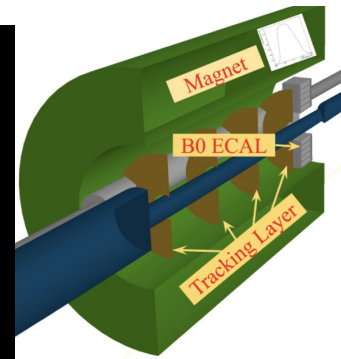


Design params: Mirror, sensor placement, gas, mirror material  
Objectives: PID performance in bins of momentum, cost



## Far Forward – B0 System

Design params: z positions of disks  
Objectives: Momentum resolution, Acceptance



# Summary and Conclusion

- Coupling the MOBO to EIC is done. Closure test 1 nearly done.
- Working on code base for a common framework for distributed optimization using PanDA and SLURM.
- EIC can be the first large-scale experiment to be realized with assistance of AI
- Ultimately, we can realize a framework that can optimize holistically a large-scale detector, and that is scalable and distributed. The Detector-2 at EIC an ideal candidate
- Exploring solutions on EIC detector Pareto front across budget values yields significant cost advantages during construction phases.
- Efficient objectives = cost-effective EIC beam time.
- This framework inherently offers broader impacts, can be adapted in various experiments and suitable for compute-intensive applications that necessitate MOO (e.g., calibrations, alignments, etc)

# Backups

# GP as a Surrogate Model

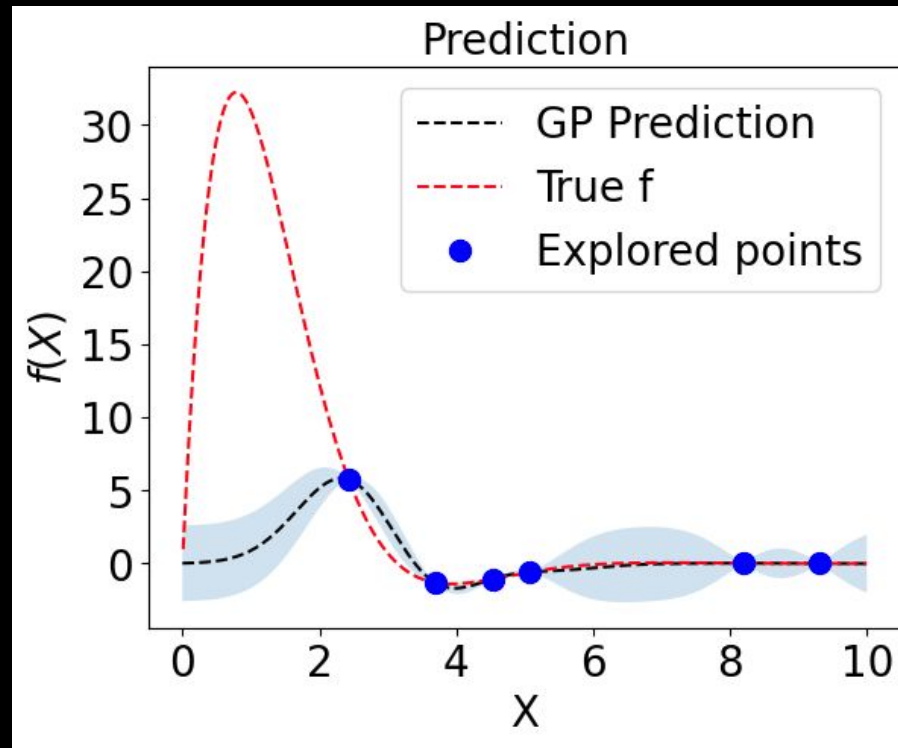
Optimization problem:

$$\underset{x}{\text{minimize}} \quad f(x)$$

$$\text{subject to} \quad g_i(x) \leq 0, \quad i = 1, \dots, m$$
$$h_j(x) = 0, \quad j = 1, \dots, p$$

**Question:** What would be the next point to explore from this?

Choose a region?

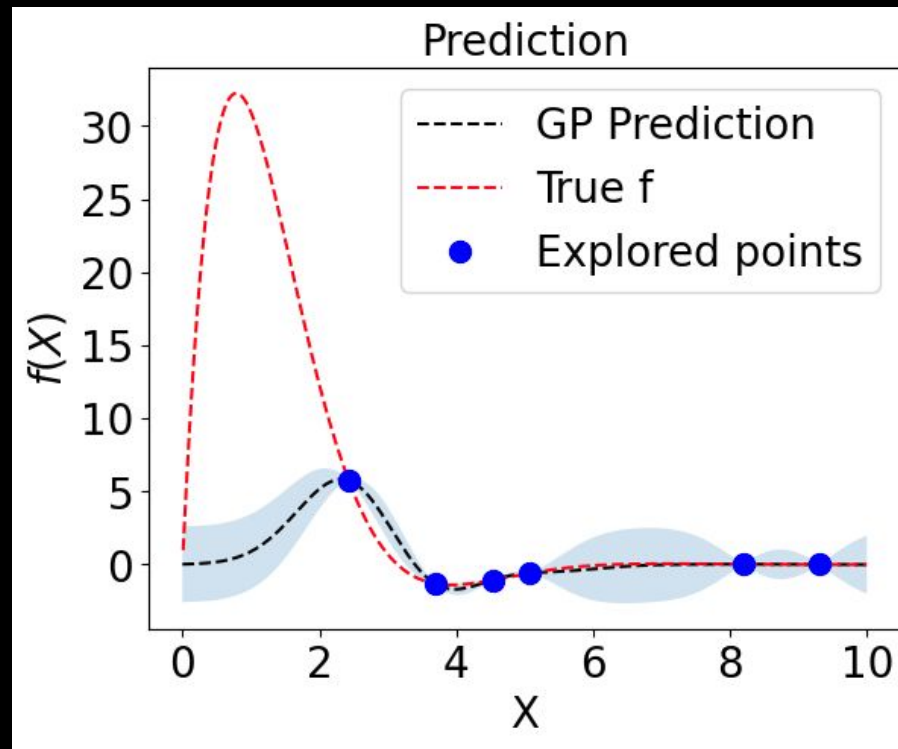


In practice we do not know the True f.

# GP as a Surrogate Model

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && g_i(x) \leq 0, \quad i = 1, \dots, m \\ & && h_j(x) = 0, \quad j = 1, \dots, p \end{aligned}$$

**The task:** To minimize. So should we even care on regions which are not minimum?

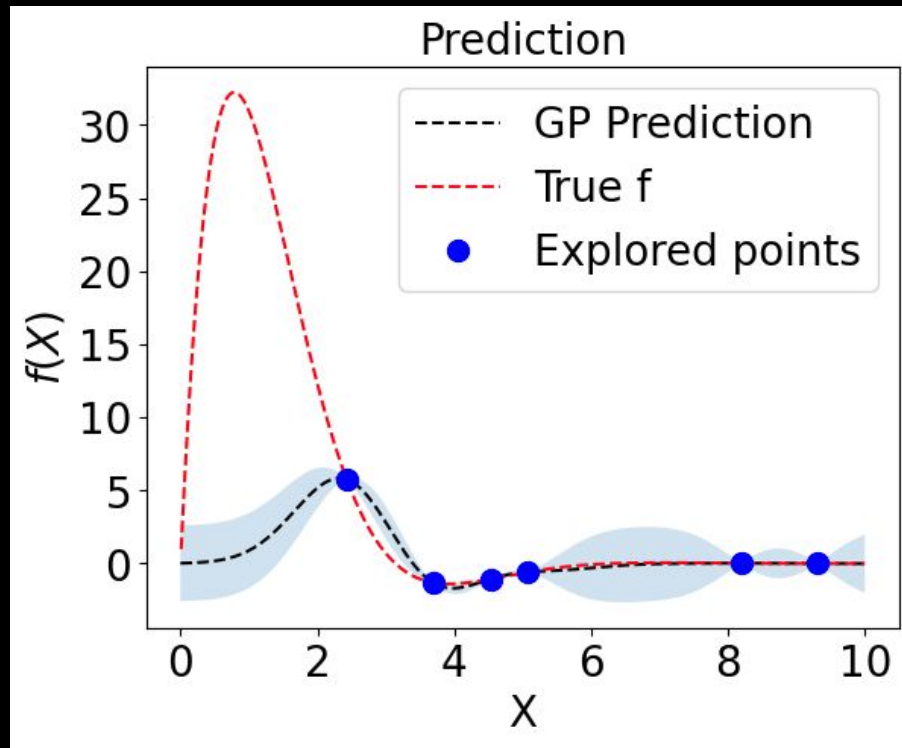


# The Acquisition function

- Define a function that scans through the search space for values of  $f(x)$  using the built GP.
- Much faster than evaluations.
- Carefully choose the next point to evaluate\*.
- Model inaccurate in region out of interest

## Widely used Acquisition functions

- Confidence Bound
- Probability of Improvement
- Expected Improvement



\*Since evaluations are supposed to be very costly

AI Assisted Detector Design for EIC



# Confidence Bound

$$\text{LCB} = -\mu(x) + \lambda\sigma(x)$$

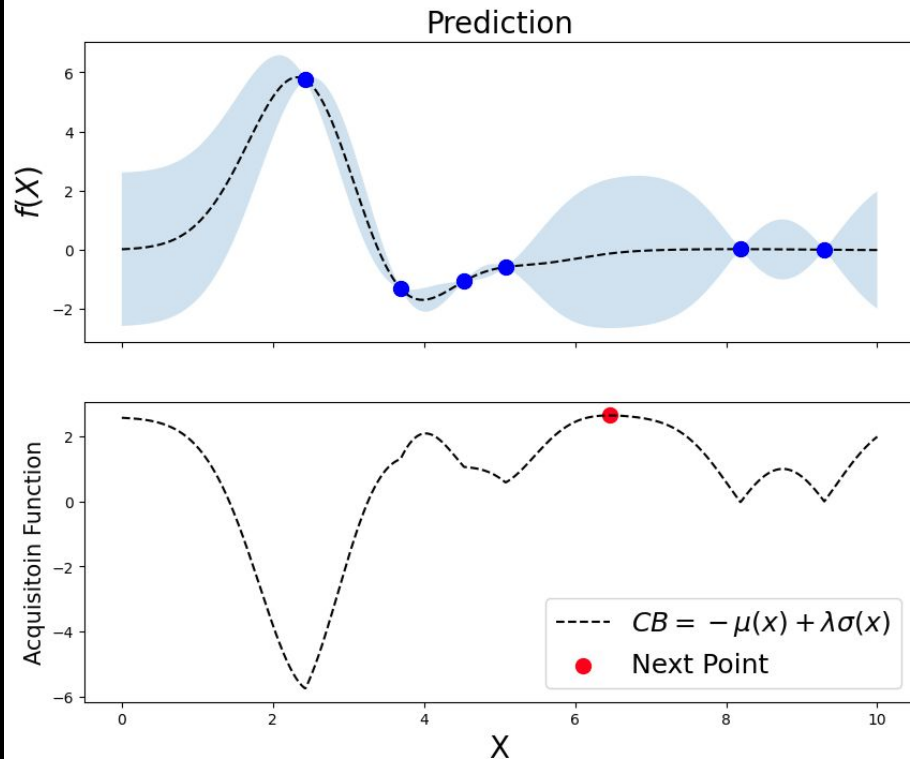
Hyper parameter  $\lambda$

Exploitation  $-\mu(x)$

Exploration  $\lambda\sigma(x)$

Can now control where the search will happen in subsequent **iterations**.

Usually, bias is towards the mean. Since it is optimization



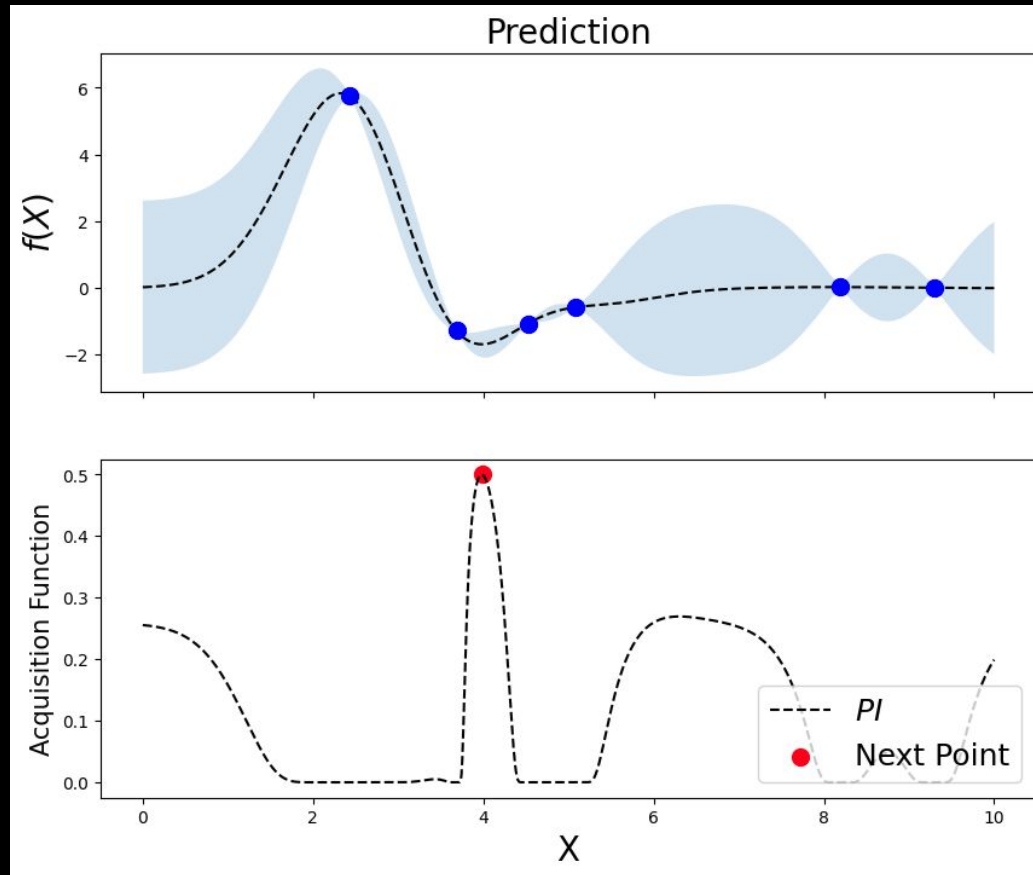
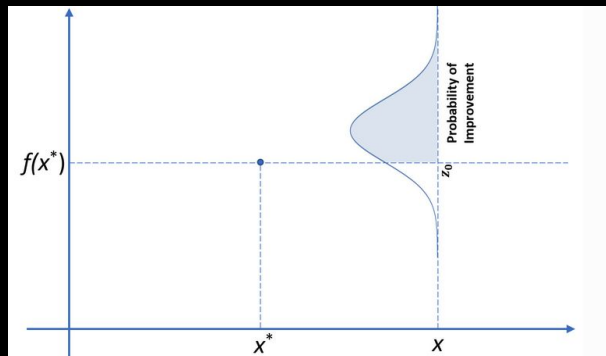
```
def UCB_Acq(input_range, gp_model, Lambda = 1.0):  
    xs = np.linspace(start = input_range[0], stop = input_range[-1], num = 1000)  
    y, y_std = gp_model.predict(xs.reshape(-1, 1), return_std = True)  
    UCB = -y + Lambda*y_std  
    ✨ return UCB
```

# Probability of Improvement

$$PI = \text{CDF}\left(\frac{\min(f(x)) - f(x^*)}{\sigma(x) + \epsilon}\right)$$

Choose the point, that has the maximum probability of improvement.

Note: NO consideration of actual value.

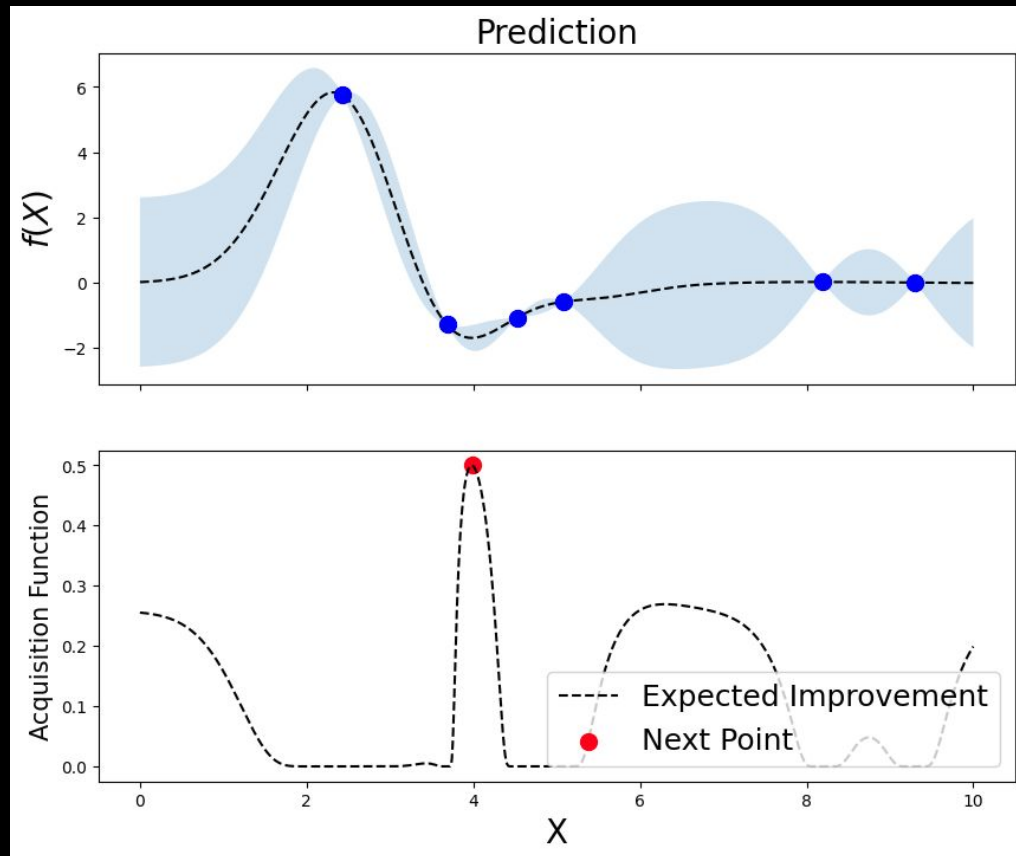


# Expected Improvement

$$EI = (\min(f(x) - f(x^*)))\text{CDF}\left(\frac{\min(f(x) - f(x^*))}{\sigma(x) + \epsilon}\right) + \sigma(x)\text{PDF}\left(\frac{\min(f(x) - f(x^*))}{\sigma(x) + \epsilon}\right)$$

Considers the Magnitude of improvement along with its probability<sup>[1]</sup>

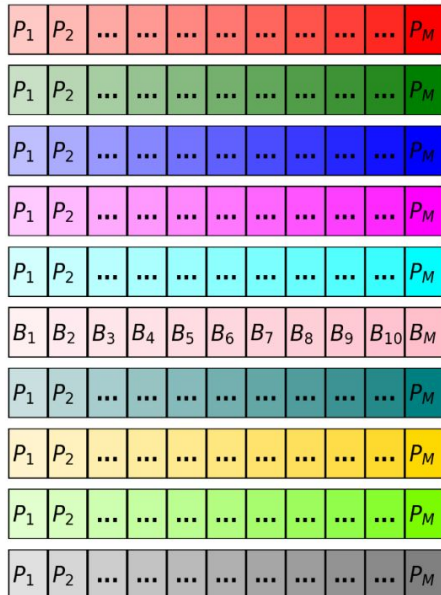
- Now, With the suggested point, Start running iterations. Choose the first q points suggested by the Acquisition function.
- Run for N iterations
- Implement early stopping criterion if necessary



# The Summary of MOGA Pipeline

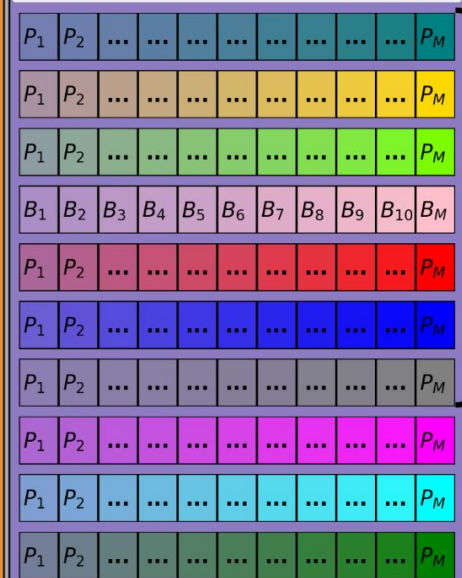
## Initial population creation (N\_pop)

Inject baseline Genes  
Faster convergence



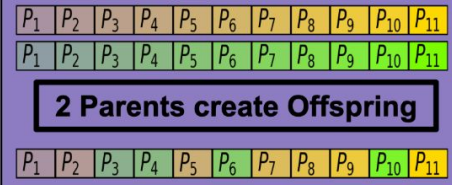
## The Evolution Cycle

### Rank & sort - NSGA2 (Objs)

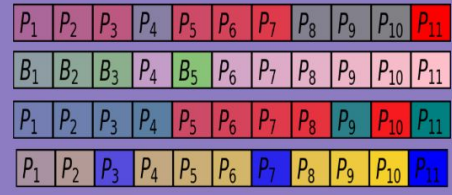


Healthy Design points

### Genetic Evolution of designs



### N\_Offsprings for next call



ePIC dd4hep Sim + eic-recon

Yields Performance of the design.  
Objectives that decide evolution

# Multi Objective Evolutionary Algorithms

- Inspired by Biological Systems.
- Semi heuristic in nature.
- Quite successful in solving MOO problems.
- Embedding constraints relatively easier



jMetalPy



## Swarm Algorithms

Ant Colony optimization  
Bees algorithm  
Particle swarm optimization  
Cuckoo search

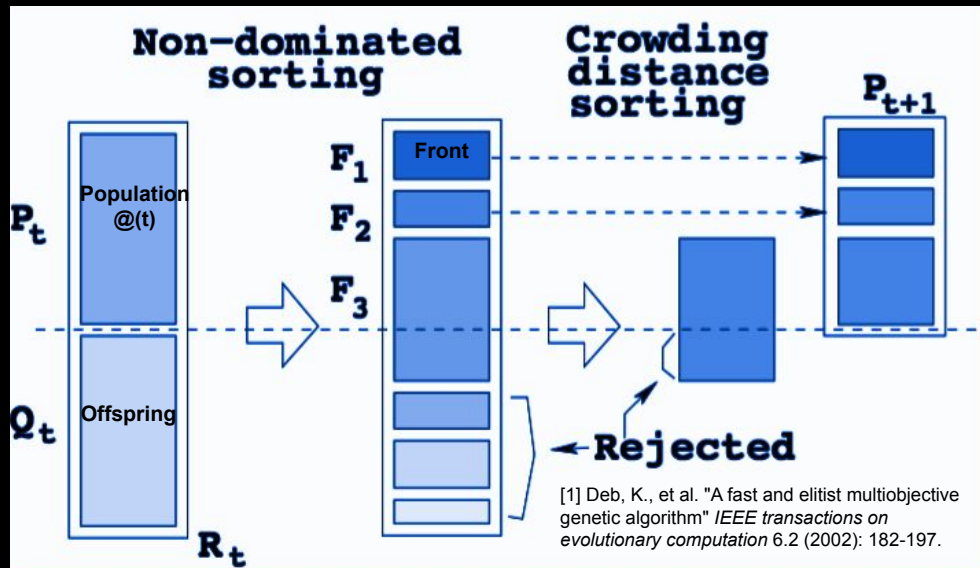
## Genetic Algorithms

Default Genetic Algorithm  
NSGA  
**NSGA-II**  
U-NSGA-III

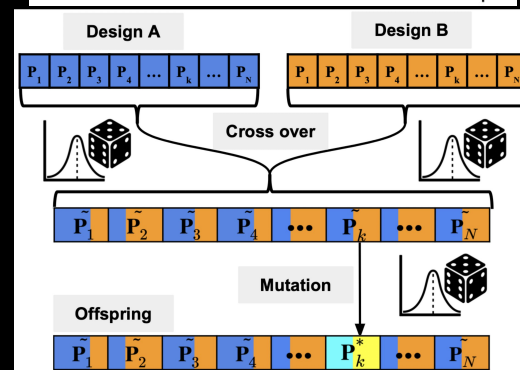
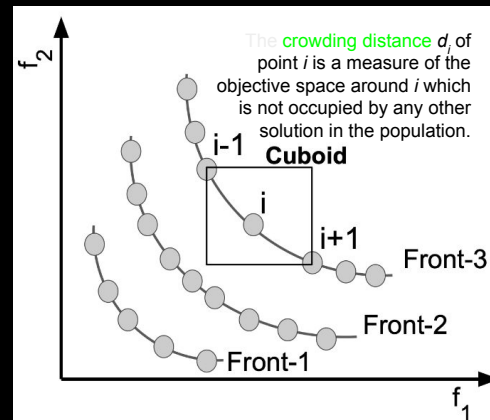
## Differential Evolution

## Cellular Automata

# Elitist Non-Dominated Sorting Genetic (NSGA)



The population  $R_t$  is classified in non-dominated fronts. Not all fronts can be accommodated in the  $N$  slots of available in the new population  $P_{t+1}$ . We use **crowding distance** to keep those points in the last front that contribute to the highest diversity.



This is to illustrate Binary Cross-over

This is one of the most popular approach (>35k citations on google scholar), characterized by:

- Use of an elitist principle
- Explicit diversity preserving mechanism
- Emphasis in non-dominated solutions

# MOEA or MOBO ?

## MOEA

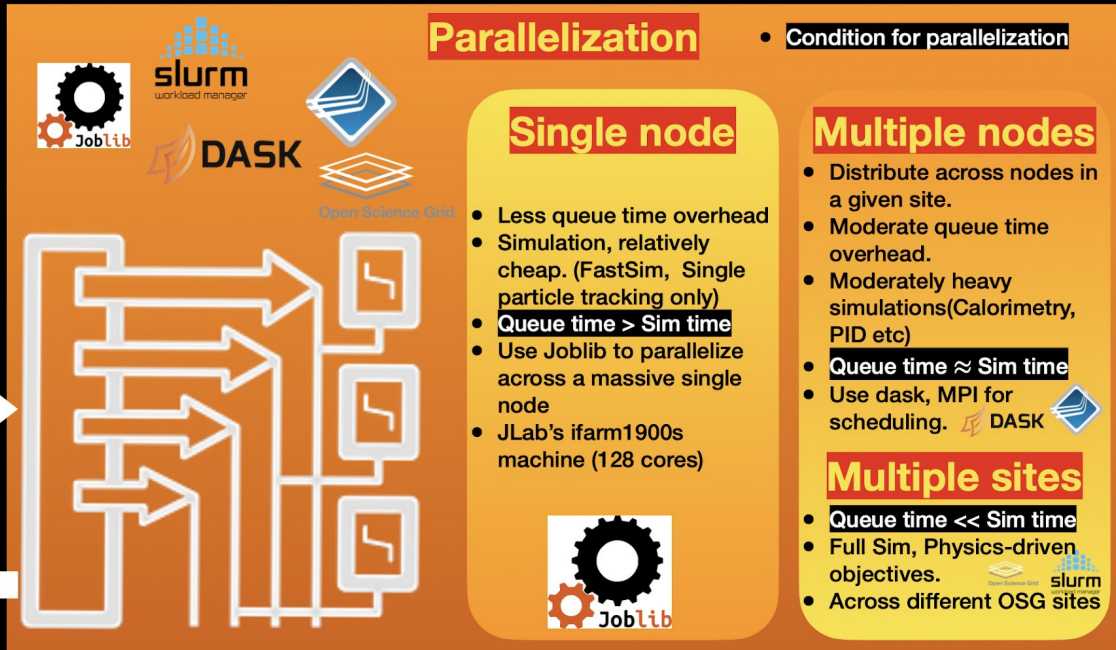
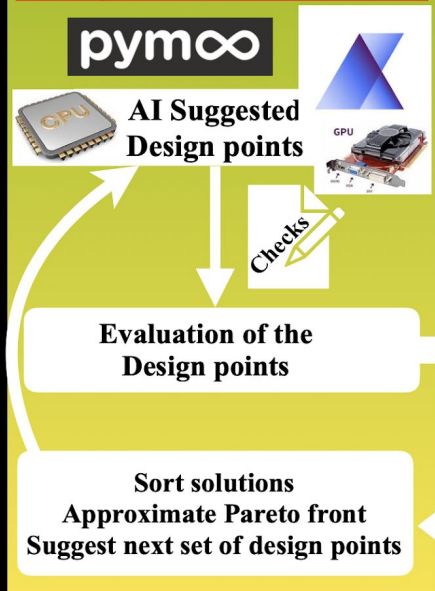
- Has been widely used for solving MOO problems
- population /off spring — diversity —
- Relatively easier to implement
- Complexity relatively easy to compute
- Ideal — Cost of computing “cheap”
- Successful with large Design and Objective parameters
- No Map : “Design” “Objectives”

## MOBO

- Has been around for a while, gaining popularity
- Sequential Strategy — global minimization
- Relatively harder to implement
- Complexity relatively easy to compute
- Ideal — simulations can be heavily parallelized
- Currently, Not recommended beyond 4-5 Objective parameters
- Can Map : “Design” “Objectives” — Fast simulator can be built

# A roadmap for scalable optimization

## AI Optimization block



Need for better visualizations  
Beyond 3D Pareto visualizations



# Far Forward Updates

<b>Problem</b>	Optimize the momentum resolution subject to the non-homogenous Magnetic field and to increase occupancy at B0 ECAL.		
<b>Objective Space = 2</b>	<b>Objective Parameter</b>	<b>Remarks</b>	
	Momentum resolution ( $p_T$ )	Momentum range of 80 - 100 GeV/c is of interest and specifically proton tracks	
	B0 ECAL acceptance	Ratio of number of tracks before 1st tracking disk to the number of showers detected by B0ECAL	
<b>Design Space = 4</b>	<b>Design Parameter</b>	<b>Range [cm]</b>	<b>Least count for variation [cm]</b>
	$Z_1$	583.0 - 630.0	1.0
	$\Delta Z_2, \Delta Z_3, \Delta Z_4$	10.0 - 40.0	1.0
<b>Constraints = 2</b>	$Z_1 + \sum_{i=2,3,4} \Delta Z_i \leq 685.5 \text{ cm}$		
	$ Z_{i+1} + Z_i  \geq 10.0 \text{ cm}$		