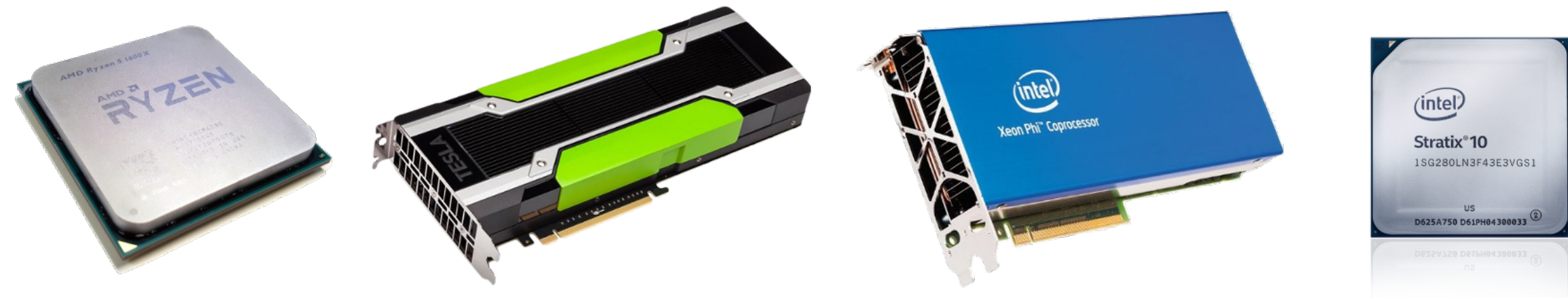# Autotuning under Tight Budget Constraints: A Transparent Design of Experiments Approach

*Pedro Bruel*, Steven Quinito Masnada, Brice Videau, Arnaud Legrand, Jean-Marc Vincent, Alfredo Goldman

## Autotuning: Optimizing Program Configurations

- ▶ How to write efficient code for each of these?
- ▶ We can use autotuning: the process of automatically finding a configuration of a program that optimizes an objective

## Strategies for Exploring Search Spaces

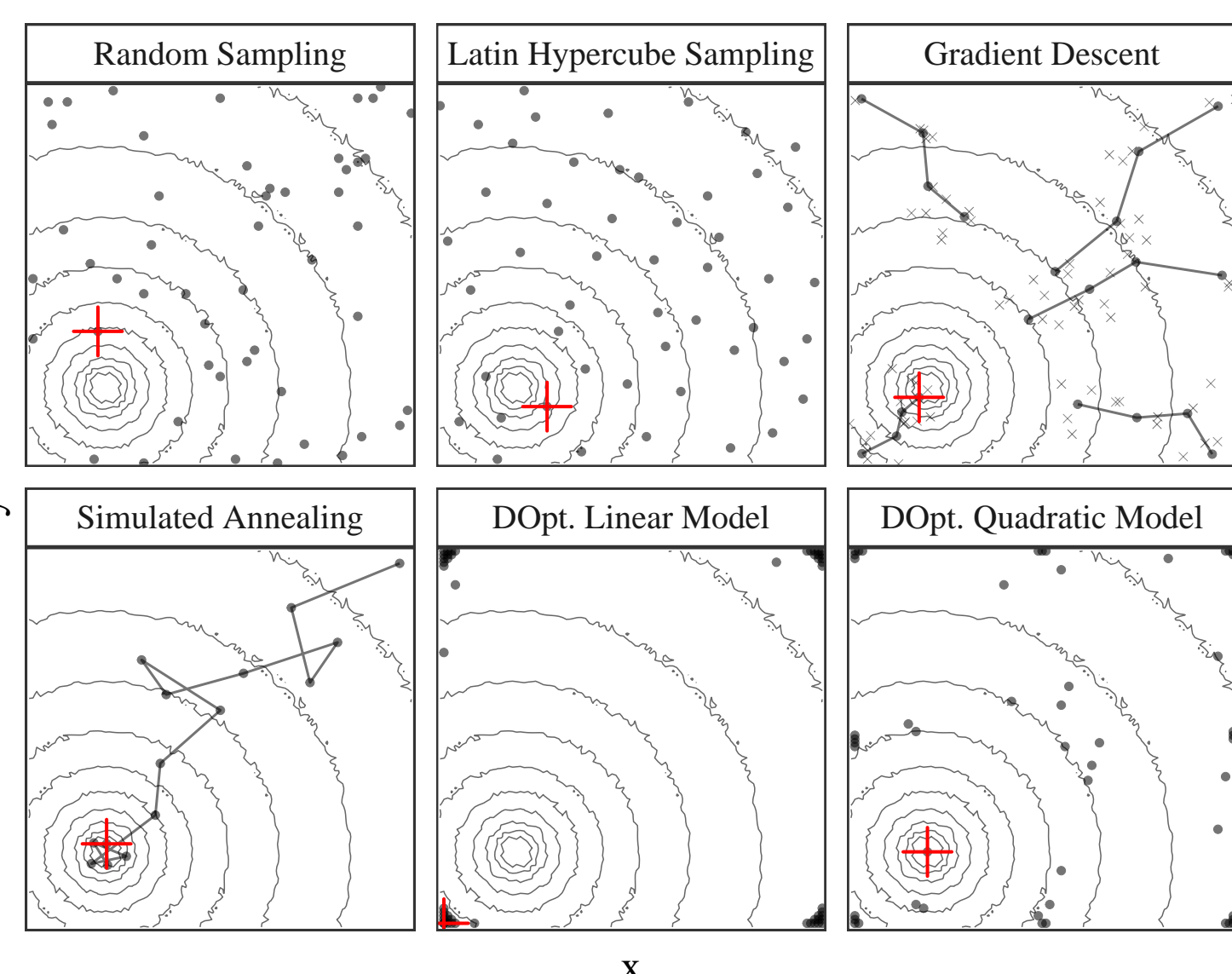| System | Domain | Approach |
|---|---|---|
| ATLAS | Dense Linear Algebra | Exhaustive |
| INSIEME | Compiler | Genetic Algorithm |
| Active Harmony | Runtime | Nelder-Mead |
| ParamILS | Domain-Agnostic | Stochastic Local Search |
| OPAL | Domain-Agnostic | Direct Search |
| OpenTuner | Domain-Agnostic | Ensemble |
| MILEPOST GCC | Compiler | Machine Learning |
| Apollo | GPU kernels | Decision Trees |

Exhaustive, Meta-Heuristics, Machine Learning

- ▶ These approaches need a large number of function evaluations, assuming seach space "smoothness", and that good solutions are reachable
- ▶ After optimizing, we learn "nothing" about the search space, and can't explain why optimizations work

## Design of Experiments: Exploration under a Budget



A Plackett-Burman design for 7 2-level factors

- ▶ Experiment results can be used to identify relevant parameters and to fit a linear regression model



- ▶ Exploration of a search space using a fixed budget of 50 points, the red "+" represents the best point found by each strategy

## A Motivating Result on a GPU Kernel
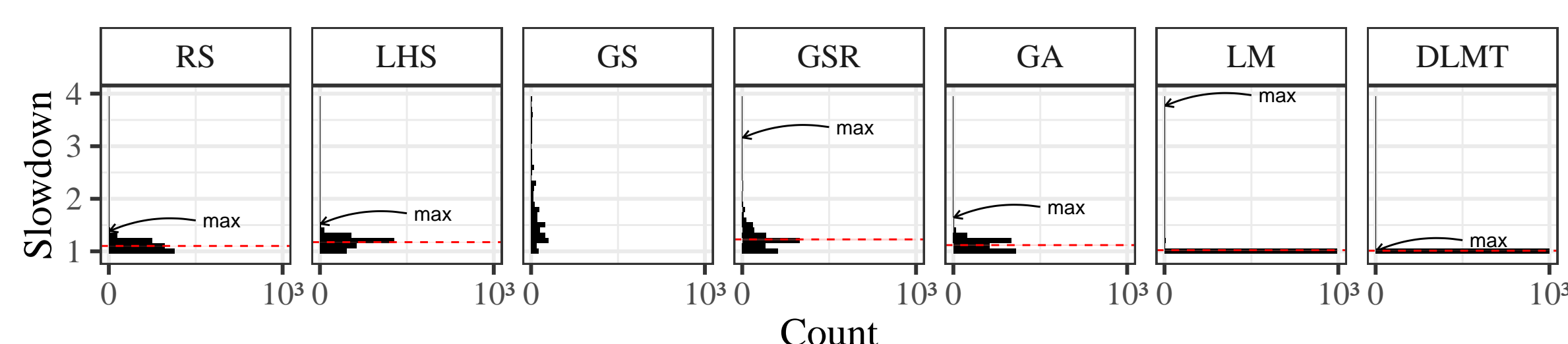
- ▶ Kernel factors:

| Factor | Levels | Short Description |
|---|---|---|
| vector_length | $2^0, \ldots, 2^4$ | Size of support arrays |
| load_overlap | *true, false* | Load overlaps in vectorization |
| temporary_size | 2, 4 | Byte size of temporary data |
| elements_number | 1, ..., 24 | Size of equal data splits |
| y_component_number | 1, ..., 6 | Loop tile size |
| threads_number | $2^5, \ldots, 2^{10}$ | Size of thread groups |
| lws_y | $2^0, \ldots, 2^{10}$ | Block size in y dimension |

- ▶ Initial performance model:

$$time\_per\_pixel \sim y\_component\_number + \frac{1}{y\_component\_number} + \\ load\_overlap + temporary\_size + \\ vector\_length + lws\_y + \frac{1}{lws\_y} + \\ elements\_number + threads\_number + \\ \frac{1}{elements\_number} + \frac{1}{threads\_number}.$$
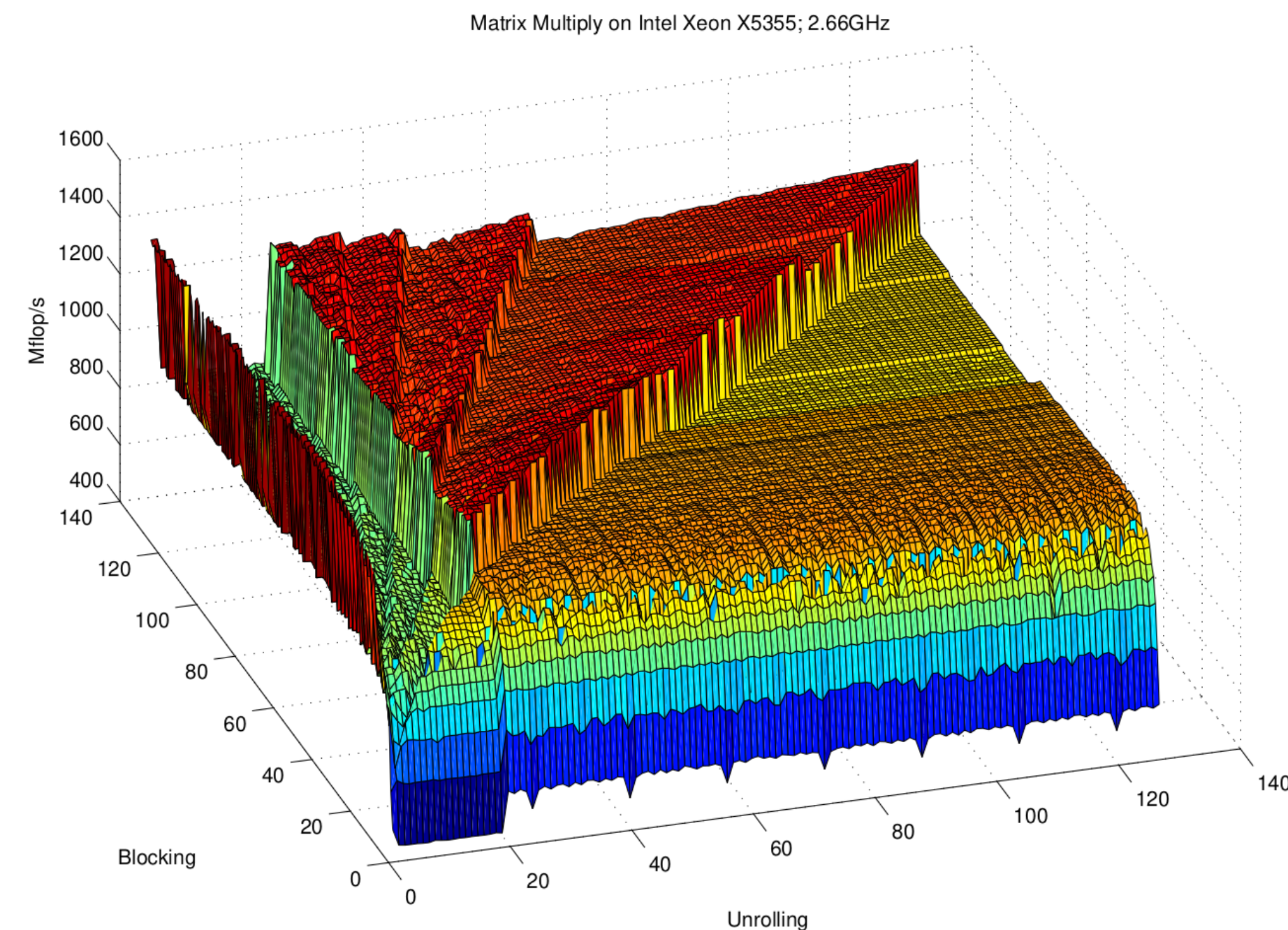
- ▶ This simple case had known valid search space and global optimum, and fixed budget

Our approach (DLMT) was always within 1% of the optimum



| RS | LHS | GS | GSR | GA | LM | DLMT |
|---|---|---|---|---|---|---|
| Random Sampling | Latin Hyper Square | Greedy Search | Greedy with Restart | Generic Algorithm | Linear Model | Our DoE Approach |

## Autotuning: Search Spaces are Hard to Explore



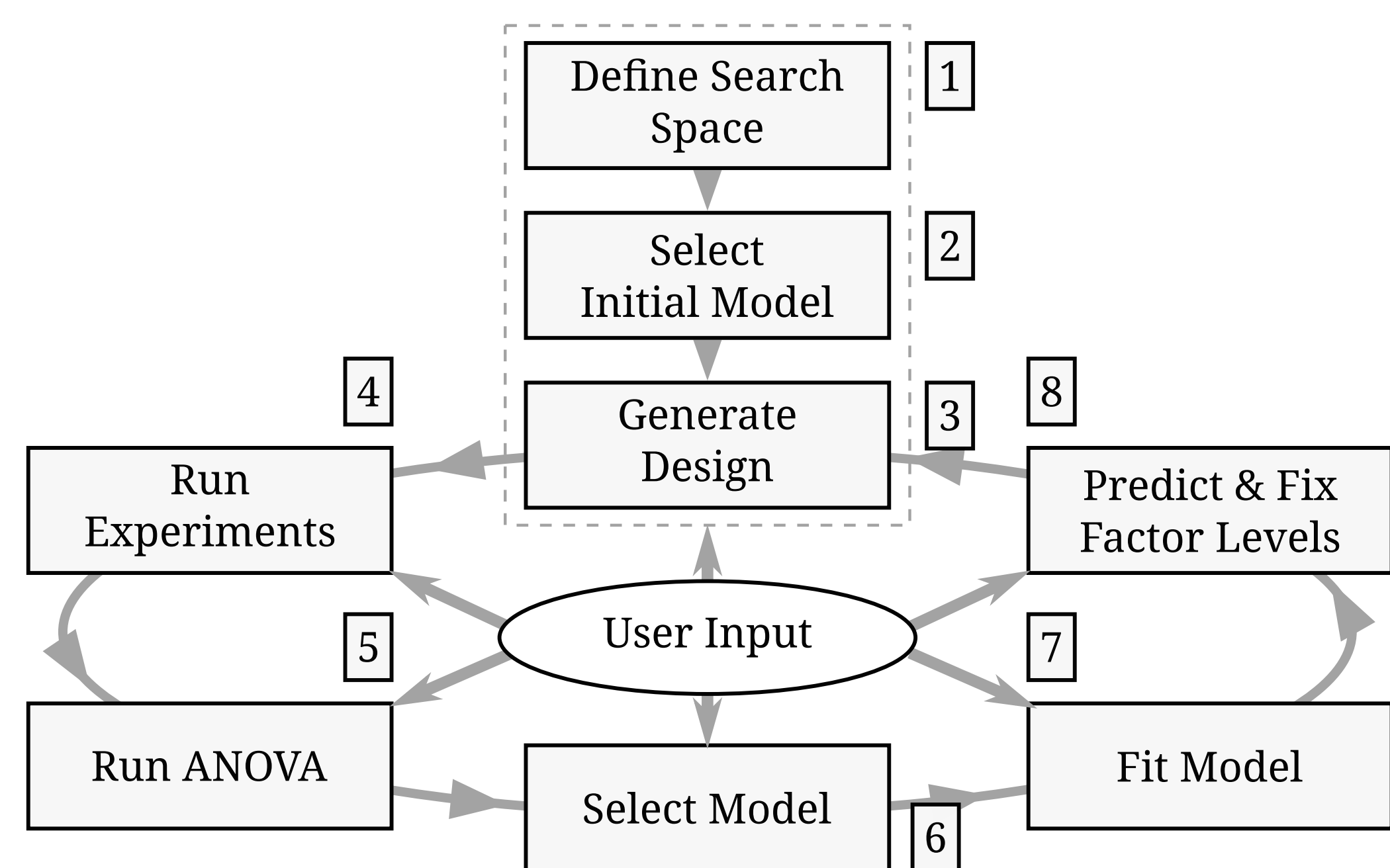Matrix Multiply on Intel Xeon X5355; 2.66GHz

Unrolling, blocking and Mflops/s for matrix multiplication
Seymour K, You H, Dongarra J. A comparison of search heuristics for empirical code optimization. InCLUSTER 2008 Oct 1 (pp. 421-429)

- ▶ Represent the effect of all possible configurations on the objectives, can be difficult to explore, with multiple local optima and undefined regions
- ▶ Main issues are exponential growth, geometry, & measurement time

## A Transparent Design of Experiments Approach



- ▶ An initial model is provided by the user (steps 1 & 2)
- ▶ Design of Experiments (DoE) guides exploration (steps 3 & 4)
- ▶ Significant factors are identified by Analysis of Variance (ANOVA) (steps 5 & 6)
- ▶ New fitted model predicts best value for significant factors (steps 7 & 8)

Transparent: factor and level selections based on ANOVA

Parsimonious: DoE decreases measurements

## Extensive Evaluation on the SPAPT Benchmark

- ▶ SPAPT is an autotuning benchmark for CPU kernels, with search space sizes between $10^7$ and $10^{36}$
- ▶ We evaluated DLMT on 11 kernels (3 shown below) using the same initial performance model, and fixed budget

Our approach (DLMT) achieved good speedups using a smaller budget, while exploring better configurations

Systèmes Répartis, Calcul Parallèle et Réseaux - POLARIS