

# MLKAPS: Machine Learning for Kernel Accuracy and Performance Studies

Mathys Elliott Jam<sup>1</sup>   Eric Petit<sup>2</sup>   Pablo Oliveira<sup>1</sup>   David Defour<sup>3</sup>   Greg Henry<sup>2</sup>   William Jalby<sup>1</sup>

<sup>1</sup>LI-PARAD

Université Versailles St-Quentin en Yvelines (UVSQ)

<sup>2</sup>Intel Corp.

<sup>3</sup>Université Perpignan Via Domitia

June 26, 2025



# **Introduction**

---

## Exemple Jouet

---

```
1 void kernel(double* data, int m, int n, int nthreads)
2 {
3     #pragma omp parallel num_threads(nthreads)
4     ...
5 }
```

L'utilisateur doit fournir tout les paramètres au noyaux. On distingue :

- **Les paramètres d'entrées (data, m, n)**
- **Les paramètres de designs (nthreads)**

# Exemple Jouet

---

```
1 void kernel_wrapped(double* data, int m, int n)
2 {
3     int nthreads = dt_nthreads(data, m, n);
4     kernel(data, m, n , nthreads);
5 }
```

**On cherche un mécanisme de décision runtime** - ici `dt_nthreads(...)` - pour trouver automatiquement les **paramètres de designs** selon les **paramètres d'entrées**.

## Exemple Jouet

---

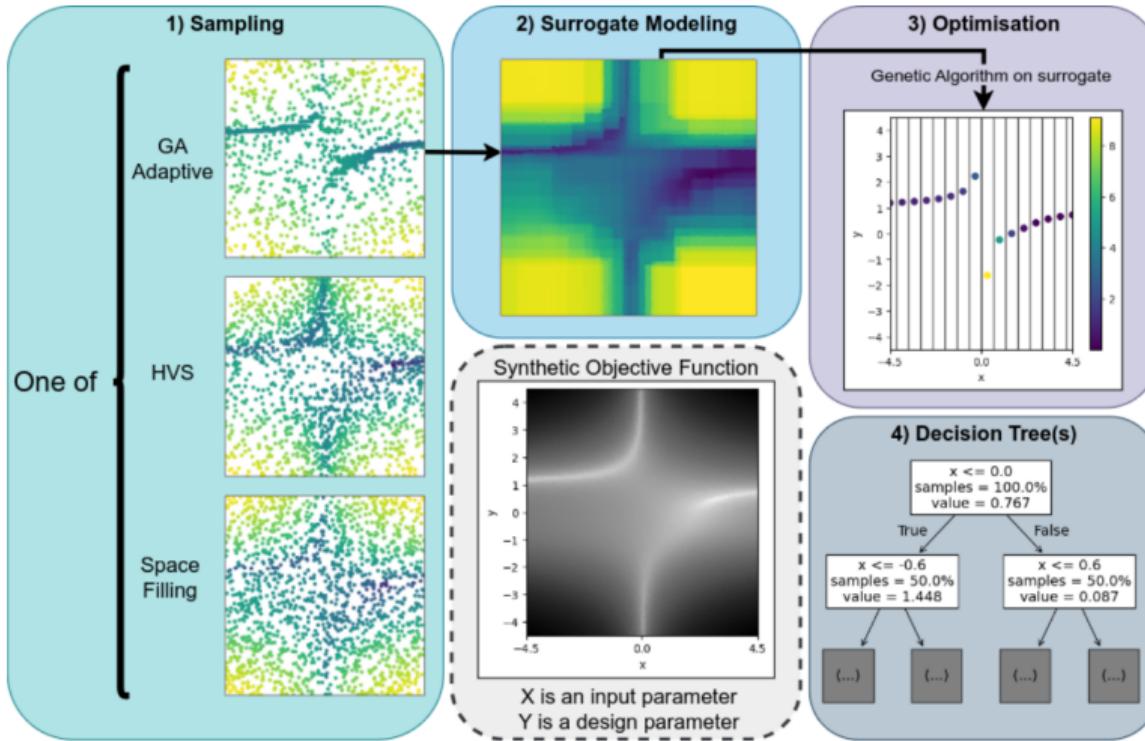
```
1 void kernel_wrapped(double* data, int m, int n)
2 {
3     int nthreads = dt_nthreads(data, m, n);
4     kernel(data, m, n , nthreads);
5 }
```

**On cherche un mécanisme de décision runtime** - ici `dt_nthreads(...)` - pour trouver automatiquement les **paramètres de designs** selon les **paramètres d'entrées**.

Soit  $O$  une fonction objective du noyau. On cherche une fonction  $DT : I \rightarrow D$  telle que, pour chaque  $i_k \in I$  :

$$DT(i_k) = d_k = \arg \min_{d \in D} O(i_k, d)$$

# Workflow de MLKAPS



## Résultats

---

# Environnement

---

	CPU	freq.	cores	threads	L1	L2	L3	RAM	Type
<b>KNM</b>	Intel Knights Mill	1.5GHz	72	288	32 KB	36 MB	—	16 GB	HBM
<b>SPR</b>	Intel Xeon Gold 6438M	2.2GHz	64	128	80 KB	2 MB	60 MB	504 GB	DDR5

On s'intéresse aux noyau dgetrf (LU) de la Intel MKL.

- 2 entrées (tailles de la matrice  $m \times n$ , avec  $m, n \in [1000, 5000]$ )
- 8 paramètres de design (nombre de threads, taille de blocs, seuil de changement de stratégie, etc.)<sup>1</sup>

Total de  $4.6 \times 10^{13}$  combinaisons de design.

---

1. Liste exacte confidentielle.

# Geometric Mean Speedup

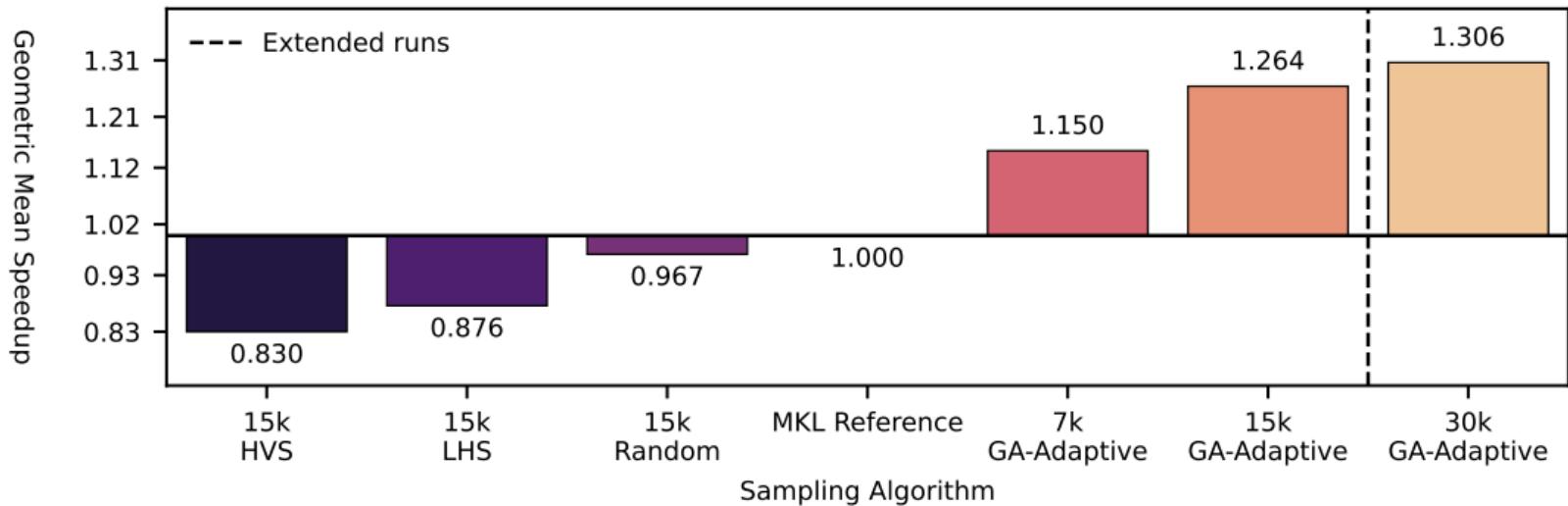


Figure – (SPR) Geomean Speedup (Higher is better) relativement aux configurations de références de la MKL sur dgetrf (LU) selon l'algorithme échantillonnage. Grille de validation de 46x46 sur  $I$ .

# Speedup sur l'espace

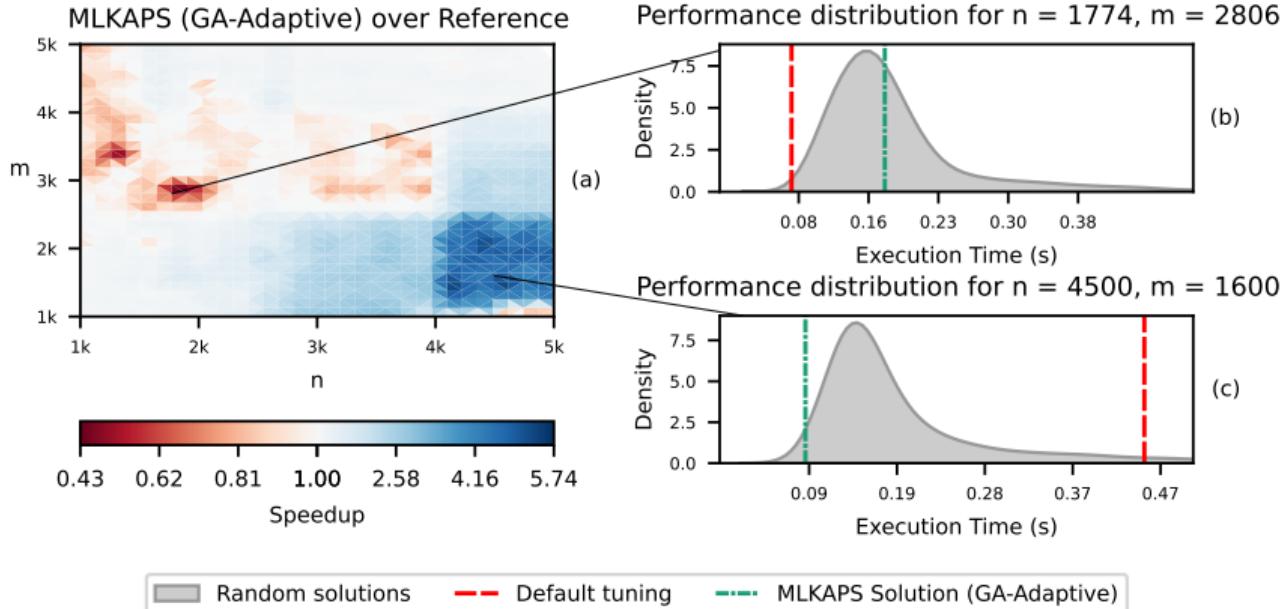


Figure – (KNM) dgetrf (LU) a) Speedup de GA-Adaptive (7k samples) comparer a la MKL.  
b/c) Distribution de la performance des configurations local.

# Comparaison a l'état de l'art

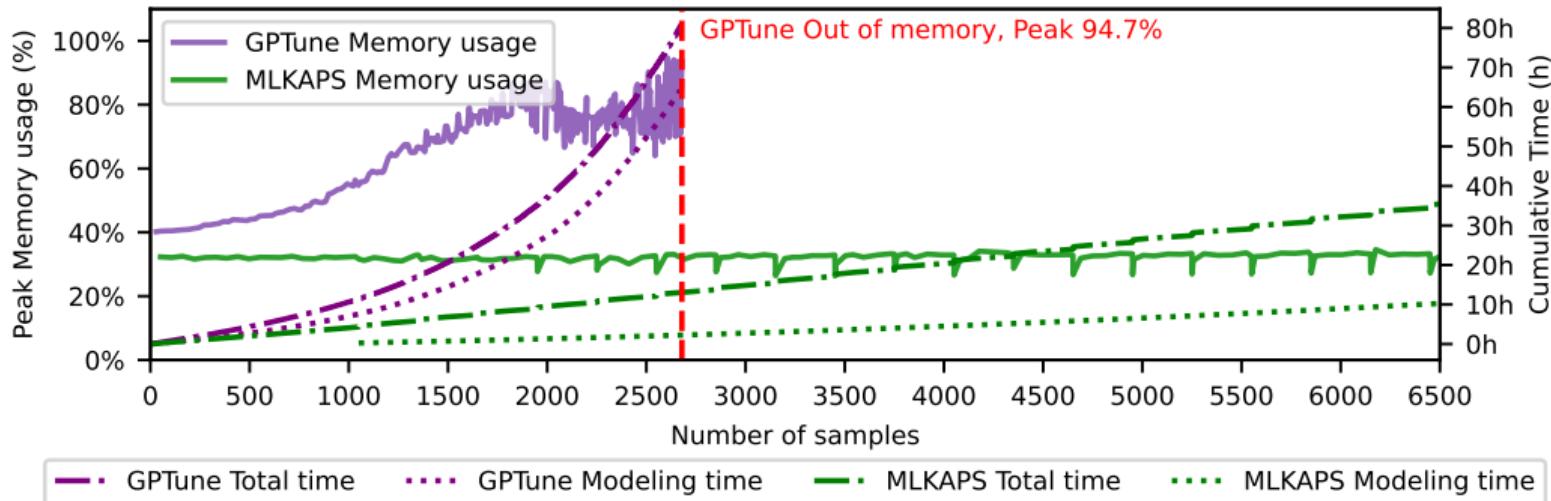


Figure – (KNM) dgetrf (LU) Évolution de la consommation mémoire et du temps d'exécution selon le nombre d'échantillons pour MLKAPS et GPTune (16 tasks).

1. GPTune : multitask learning for autotuning exascale applications, Yang Liu, Wissam M Sid-Lakhdar, Osni Marques, Xinran Zhu, Chang Meng, James W Demmel, and Xiaoye S Li.

# Conclusions

---

Soumission ACM Taco under review, disponible sur HAL<sup>2</sup> et arXiv<sup>3</sup>

- Projet MLKAPS débuté par Intel en 2022, maintenant dans le domaine public <https://github.com/MLCGO/MLKAPS>
- Les arbres de décisions de MLKAPS pour QR et LU sont dans la Intel MKL depuis la release 2025.2
- Plus de résultats dans le papier.

Sinon : [mathysjam.org](http://mathysjam.org)

- 
2. <https://hal.science/hal-04851397v1>
  3. <https://arxiv.org/abs/2501.05811>