

# Performance-portability and task parallelism for Lattice Field Theory

Bartosz Kostrzewa<sup>\*</sup>, Stefan Krieg<sup>†</sup>, Estela Suarez<sup>‡</sup>

<sup>\*</sup>High Performance Computing & Analytics Lab

<sup>†</sup>Helmholtz Zentrum für Strahlen und Kernphysik

<sup>‡</sup>Institut für Informatik

Rheinische Friedrich-Wilhelms-Universität Bonn

<sup>†‡</sup>Juelich Supercomputing Center, Forschungszentrum Juelich

Colour meets Flavour – Innovative Algorithms Lab Kick-off meeting  
2024.03.20



**<NUMERIQS>**

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).
  - ▶ Strong scalability: most efficient algorithms don't scale well (maybe they don't need to).

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).
  - ▶ Strong scalability: most efficient algorithms don't scale well (maybe they don't need to).
  - ▶ **Complex observables** lead to hard-to-debug nested loops with horrible inter-dependencies and low re-use potential.

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).
  - ▶ Strong scalability: most efficient algorithms don't scale well (maybe they don't need to).
  - ▶ **Complex observables** lead to hard-to-debug nested loops with horrible inter-dependencies and low re-use potential.
  - ▶ **Lack of task parallelism**: lots of untapped resources which are inaccessible with current LFT frameworks.

# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).
  - ▶ Strong scalability: most efficient algorithms don't scale well (maybe they don't need to).
  - ▶ **Complex observables** lead to hard-to-debug nested loops with horrible inter-dependencies and low re-use potential.
  - ▶ **Lack of task parallelism**: lots of untapped resources which are inaccessible with current LFT frameworks.
  - ▶ **Storage**: we cannot store millions of gauge configurations → might need on-the-fly calculations of **everything**.



# Technical Challenges in Lattice Field Theory

⟨NUMERICS⟩

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).
  - ▶ Strong scalability: most efficient algorithms don't scale well (maybe they don't need to).
  - ▶ **Complex observables** lead to hard-to-debug nested loops with horrible inter-dependencies and low re-use potential.
  - ▶ **Lack of task parallelism**: lots of untapped resources which are inaccessible with current LFT frameworks.
  - ▶ **Storage**: we cannot store millions of gauge configurations → might need on-the-fly calculations of **everything**.
  - ▶ Programmer productivity: need high-performance backends with easy to use frontends for students (and for ourselves).

# Technical Challenges in Lattice Field Theory

(NUMERICS)

- Assume that sampling problems in lattice field theory can be tempered or overcome:
  - ▶ critical slowing down
  - ▶ ergodicity issues in the presence of multiple phases
  - ▶ volume scaling of generative flow models
  - ▶ multi-level sampling
- What are the major technical challenges facing moonshot-type calculations?
  - ▶ **HPC heterogeneity**: need write-once run-anywhere (as far as possible).
  - ▶ Strong scalability: most efficient algorithms don't scale well (maybe they don't need to).
  - ▶ **Complex observables** lead to hard-to-debug nested loops with horrible inter-dependencies and low re-use potential.
  - ▶ **Lack of task parallelism**: lots of untapped resources which are inaccessible with current LFT frameworks.
  - ▶ **Storage**: we cannot store millions of gauge configurations → might need on-the-fly calculations of **everything**.
  - ▶ Programmer productivity: need high-performance backends with easy to use frontends for students (and for ourselves).
  - ▶ Physics students: not enough training in software engineering / HPC / modern C++ (for this kind of work)

# HPC heterogeneity

(NUMERICS)

## Need to target different architectures



- Situation bound to get worse with future specialized hardware & Modular Supercomputing Architecture
  - ▶ multi-purpose chiplet designs (EPI)
  - ▶ in-memory computing
  - ▶ neuromorphic devices
  - ▶ tensor & stencil accelerators, large FPGAs

# HPC heterogeneity

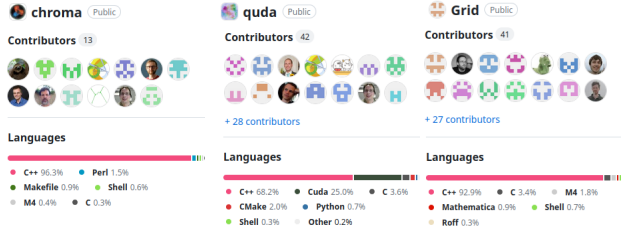
(NUMERICS)

## Need to target different architectures



- Situation bound to get worse with future specialized hardware & Modular Supercomputing Architecture
  - ▶ multi-purpose chiplet designs (EPI)
  - ▶ in-memory computing
  - ▶ neuromorphic devices
  - ▶ tensor & stencil accelerators, large FPGAs

## Existing libraries with some performance-portability



# HPC heterogeneity

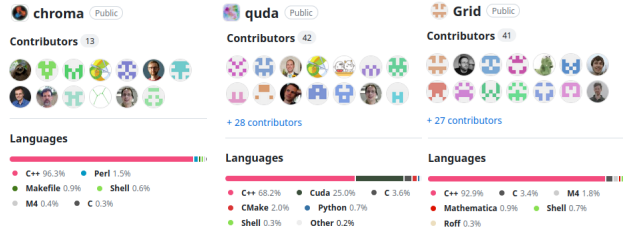
(NUMERiQS)

## Need to target different architectures



- Situation bound to get worse with future specialized hardware & Modular Supercomputing Architecture
  - ▶ multi-purpose chiplet designs (EPI)
  - ▶ in-memory computing
  - ▶ neuromorphic devices
  - ▶ tensor & stencil accelerators, large FPGAs

## Existing libraries with some performance-portability



## Future directions explored in NuMeriQS B02

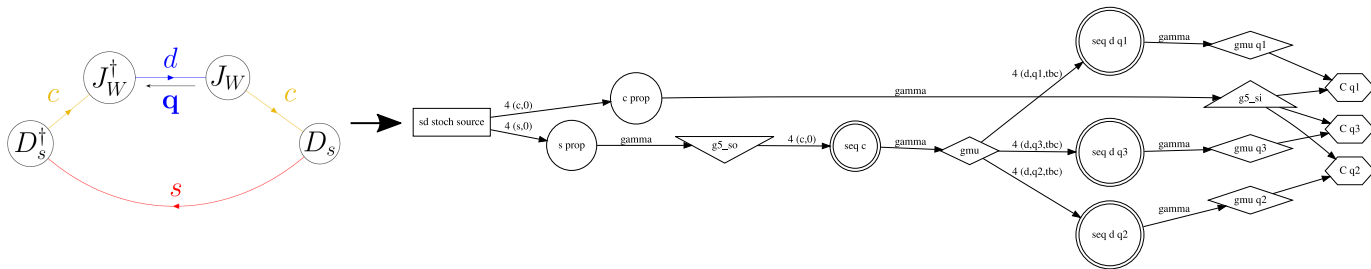
LQFT base layer: abstractions for memory space, memory access patterns and execution space, parallel patterns, asynchronicity via futures, etc.

- `std::execution`
- `std::mdspan`
- `std::async`, `std::future`
- 

# Complex observables and nested loops

(NUMERICS)

- Simple observable: two current insertions with momentum transfer between meson initial and final states
  - ▶ lots of unused parallelism at the inversion / contraction stage (different flavours, different momenta, different Dirac structures)

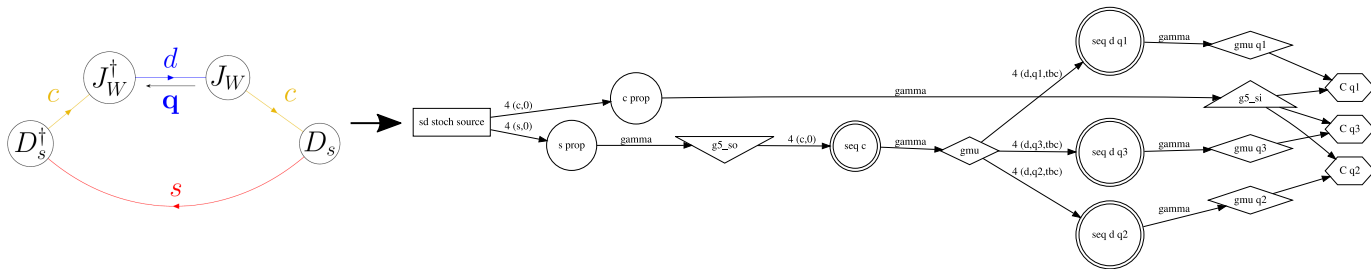


- How to exploit this without writing really complicated nested loops with tons of conditionals?

# Complex observables and nested loops

(NUMERICS)

- Simple observable: two current insertions with momentum transfer between meson initial and final states
  - ▶ lots of unused parallelism at the inversion / contraction stage (different flavours, different momenta, different Dirac structures)

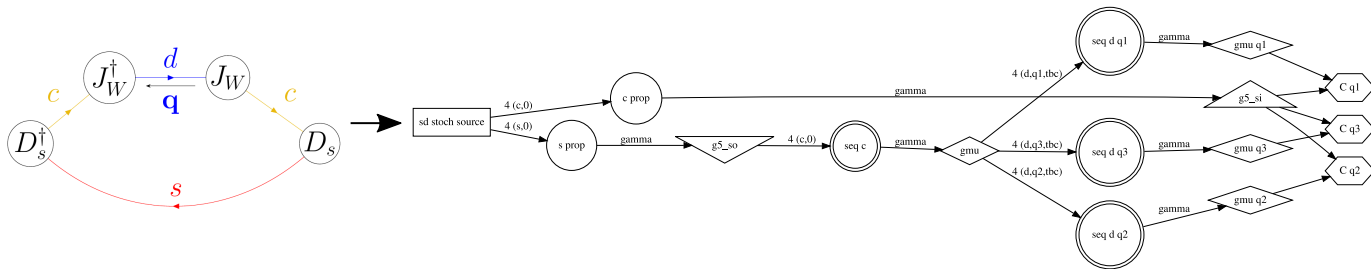


- How to exploit this without writing really complicated nested loops with tons of conditionals?
  - ▶ **Collections of observables as forests of directed acyclic graphs**

# Complex observables and nested loops

(NUMERICS)

- Simple observable: two current insertions with momentum transfer between meson initial and final states
  - ▶ lots of unused parallelism at the inversion / contraction stage (different flavours, different momenta, different Dirac structures)



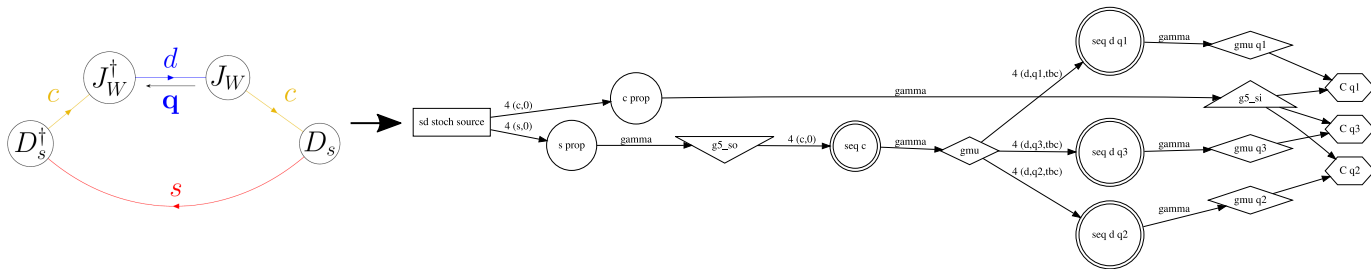
- How to exploit this without writing really complicated nested loops with tons of conditionals?
  - ▶ **Collections of observables as forests of directed acyclic graphs**
  - ▶ Edges carry computational cost



# Complex observables and nested loops

(NUMERICS)

- Simple observable: two current insertions with momentum transfer between meson initial and final states
  - ▶ lots of unused parallelism at the inversion / contraction stage (different flavours, different momenta, different Dirac structures)

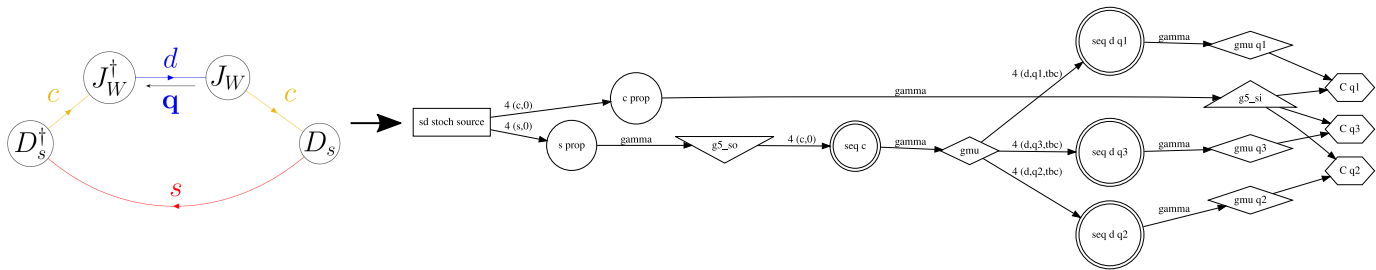


- How to exploit this without writing really complicated nested loops with tons of conditionals?
  - ▶ **Collections of observables as forests of directed acyclic graphs**
  - ▶ Edges carry computational cost
  - ▶ Vertices (may) carry computational and memory cost

# Complex observables and nested loops

(NUMERICS)

- Simple observable: two current insertions with momentum transfer between meson initial and final states
  - ▶ lots of unused parallelism at the inversion / contraction stage (different flavours, different momenta, different Dirac structures)



- How to exploit this without writing really complicated nested loops with tons of conditionals?

## ▶ Collections of observables as forests of directed acyclic graphs

- ▶ Edges carry computational cost
- ▶ Vertices (may) carry computational and memory cost

⇒ Optimise calculation under machine constraints

- ▶ all dependencies taken care of automatically, no more nested loops, just descend down the graph hierarchy
- ▶ can trade memory or storage for computation or vice-versa

# Task parallelism

(NUMERQS)

```
Every 2.0s: nvidia-smi                               mel2170: Tue Mar 19 15:34:59 2024
Tue Mar 19 15:34:59 2024
-----
| NVIDIA-SMI 535.104.12                Driver Version: 535.104.12   CUDA Version: 12.2          |
|-----|-----|-----|-----|-----|-----|-----|-----|
| GPU  Name                            Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf                      Pwr:Usage/Cap|  Bus-ID        Memory-Usage| GPU-Util  Compute M. | | | | | |
|---|---|---|---|---|---|---|---|
|  0   NVIDIA A100-SXM4-40GB            On                 | 00000000:03:00.0 Off  |      0          0          |
| N/A   52C    P0                       244W / 400W    | 5884MiB / 40960MiB |    91%    Default |
|                                     |                 |                 |                 |                 |
|                                     |                 |                 |                 |                 |
|  1   NVIDIA A100-SXM4-40GB            On                 | 00000000:44:00.0 Off  |      0          0          |
| N/A   52C    P0                       264W / 400W    | 5884MiB / 40960MiB |    90%    Default |
|                                     |                 |                 |                 |                 |
|                                     |                 |                 |                 |                 |
|  2   NVIDIA A100-SXM4-40GB            On                 | 00000000:84:00.0 Off  |      0          0          |
| N/A   51C    P0                       166W / 400W    | 5884MiB / 40960MiB |    90%    Default |
|                                     |                 |                 |                 |                 |
|                                     |                 |                 |                 |                 |
|  3   NVIDIA A100-SXM4-40GB            On                 | 00000000:C4:00.0 Off  |      0          0          |
| N/A   50C    P0                       150W / 400W    | 5884MiB / 40960MiB |    90%    Default |
|                                     |                 |                 |                 |                 |
|                                     |                 |                 |                 |                 |
|-----|-----|-----|-----|-----|-----|-----|-----|
| Processes:                               |
| GPU  GI    CI          PID  Type  Process name                        GPU Memory |
| ID   ID    ID          |     |     |                                     Usage     |
|-----|-----|-----|-----|-----|-----|-----|-----|
|  0   N/A  N/A         18516 C    ...6603758-quda-develop-dfef2c0/hmc_tm 5700MiB |
|  1   N/A  N/A         18518 C    ...6603758-quda-develop-dfef2c0/hmc_tm 5700MiB |
|  2   N/A  N/A         18515 C    ...6603758-quda-develop-dfef2c0/hmc_tm 5700MiB |
|  3   N/A  N/A         18517 C    ...6603758-quda-develop-dfef2c0/hmc_tm 5700MiB |
|-----|-----|-----|-----|-----|-----|-----|-----|
```

- Lots of free resources when running at strong-scaling limit.

HMC run on 4×A100 nodes at strong scaling limit







# Task parallelism

(NUMEROS)

```
Every 2.0s: nvidia-smi                               mel2170: Tue Mar 19 15:34:59 2024
Tue Mar 19 15:34:59 2024
-----
| NVIDIA-SMI 535.104.12                Driver Version: 535.104.12   CUDA Version: 12.2          |
|-----|-----|-----|-----|-----|-----|-----|-----|
| GPU  Name                            Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf                      Pwr:Usage/Cap|  Memory-Usage | GPU-Util  Compute M. | | | | | |
|---|---|---|---|---|---|---|---|
|  0   NVIDIA A100-SXM4-40GB           On          | 00000000:03:00.0 Off  |      0         | 91%      Default |
| N/A   52C   P0                        244W / 400W | 5884MiB / 40960MiB |          Disabled        |
|-----|-----|-----|-----|-----|-----|-----|-----|
|  1   NVIDIA A100-SXM4-40GB           On          | 00000000:44:00.0 Off  |      0         | 90%      Default |
| N/A   52C   P0                        264W / 400W | 5884MiB / 40960MiB |          Disabled        |
|-----|-----|-----|-----|-----|-----|-----|-----|
|  2   NVIDIA A100-SXM4-40GB           On          | 00000000:84:00.0 Off  |      0         | 90%      Default |
| N/A   51C   P0                        166W / 400W | 5884MiB / 40960MiB |          Disabled        |
|-----|-----|-----|-----|-----|-----|-----|-----|
|  3   NVIDIA A100-SXM4-40GB           On          | 00000000:C4:00.0 Off  |      0         | 90%      Default |
| N/A   50C   P0                        150W / 400W | 5884MiB / 40960MiB |          Disabled        |
|-----|-----|-----|-----|-----|-----|-----|-----|
| Processes:                               |
| GPU   GI    CI          PID    Type   Process name                               | GPU Memory Usage |
|-----|-----|-----|-----|-----|-----|-----|-----|
|  0   N/A   N/A         18516   C      ...6603758-quda-develop-dfef2c0/hmc_tm   5700MiB |
|  1   N/A   N/A         18518   C      ...6603758-quda-develop-dfef2c0/hmc_tm   5700MiB |
|  2   N/A   N/A         18515   C      ...6603758-quda-develop-dfef2c0/hmc_tm   5700MiB |
|  3   N/A   N/A         18517   C      ...6603758-quda-develop-dfef2c0/hmc_tm   5700MiB |
|-----|-----|-----|-----|-----|-----|-----|-----|
```

- Lots of free resources when running at strong-scaling limit.
  - GPU and CPU idle at different times.
    - ▶ Wouldn't it be neat to **fully** use the CPU on a Grace-Hopper node?
  - Would like to exploit this to run calculations in idle gaps, additional GPU streams.
- ⇒ Need suitable programming model to express task parallelism, work stealing, etc.
- ! Must be rather coarse-grained, existing asynchronous many task systems (likely) not suitable.

HMC run on 4×A100 nodes at strong scaling limit

# Summary

(NUMERICS)

- Assumption: sampling problems in lattice QFT will be tempered sufficiently
    - ▶ produce millions of large volume gauge configurations
    - ▶ impossible to store, need to calculate (almost) **everything** on the fly
- ⇒ Need: task parallelism, DAG representation of observables and dependencies, performance-portable kernels

Issue	Performance-portability	Task parallelism	loops → DAGs	Productivity layer
HPC heterogeneity	✓	–	–	–
Strong scaling	–	✓	✓	–
Wasted resources	✓	✓	✓	–
Complex observables	–	✓	✓	maybe
Storage	–	✓	✓	–
use ML frameworks	–	–	–	✓
Overwhelmed students	✗	✗	✗	maybe