

# GPU-Accelerated Particle-in-cell Code on Minsky

IWOPH17, ISC, Frankfurt a. M.

# Outline

## About

About JSC

About Supercomputers

## JuSPIC

Program Description

Steps

## Acceleration for GPUs

OpenACC

CUDA Fortran

Data Layout Analysis

Data Layout Conversion

## Performance Modelling

Effective Bandwidth

Clock Rates

## Conclusions & Outlook

## Contributions *TL;DR*

- PiC Code to GPU (partly)
- OpenACC, CUDA Fortran
- Data layout benchmarks on Minsky (POWER8NVL, P100)
- Peculiarities with PGI compiler on POWER
- Performance Model

# Jülich Supercomputing Centre

*Part of Forschungszentrum Jülich*

- Forschungszentrum Jülich
  - One of Europe's largest research centers ( $\approx 6000$  employees)
  - Energy, environmental sciences, health, information technology
- Jülich Supercomputing Centre
  - Two Top 500 supercomputers (JUQUEEN: #21, JURECA: #80)
  - NVIDIA Application Lab
  - POWER Acceleration and Design Centre



JÜLICH  
APPLICATION LAB  
NVIDIA

# Supercomputers Involved

**JURON**

**JURECA**

**JUHYDRA**



## JURON

- Human Brain Project prototype
- 18 nodes with IBM POWER8NVL CPUs ( $2 \times 10$  cores)
- Per Node: 4 NVIDIA Tesla P100 cards, connected via NVLink.
- GPU: 0.38 PFLOP/s peak performance
- NVME

## JURECA

## JUHYDRA

# Supercomputers Involved



## JURON

- Human Brain Project prototype
- 18 nodes with IBM POWER8NVL CPUs ( $2 \times 10$  cores)
- Per Node: 4 NVIDIA Tesla P100 cards, connected via NVLink.
- GPU: 0.38 PFLOP/s peak performance
- NVME



## JURECA

- General-purpose supercomputer
- 1872 nodes with Intel Xeon E5 CPUs ( $2 \times 12$  cores)
- 75 nodes with 2 NVIDIA Tesla K80 cards
- 1.8 (CPU) + 0.44 (GPU) PFLOP/s peak performance (#70)
- EDR InfiniBand

## JUHYDRA

# Supercomputers Involved



## JURON

- Human Brain Project prototype
- 18 nodes with IBM POWER8NVL CPUs ( $2 \times 10$  cores)
- Per Node: 4 NVIDIA Tesla P100 cards, connected via NVLink.
- GPU: 0.38 PFLOP/s peak performance
- NVME



## JURECA

- General-purpose supercomputer
- 1872 nodes with Intel Xeon E5 CPUs ( $2 \times 12$  cores)
- 75 nodes with 2 NVIDIA Tesla K80 cards
- 1.8 (CPU) + 0.44 (GPU) PFLOP/s peak performance (#70)
- EDR InfiniBand



## JUHYDRA

- GPU prototyping machine
- 1 node with Intel Xeon E5 CPU ( $2 \times 8$  cores)
- NVIDIA  $2 \times$  Tesla K20,  $2 \times$  Tesla K40 cards
- No batch system



## JURON

- Human Brain Project prototype
- 18 nodes with IBM POWER8NVL CPUs ( $2 \times 10$  cores)
- Per Node: 4 NVIDIA Tesla P100 cards, connected via NVLink.
- GPU: 0.38 PFLOP/s peak performance
- NVME



## JURECA

- General-purpose supercomputer
- 1872 nodes with Intel Xeon E5 CPUs ( $2 \times 12$  cores)
- 75 nodes with 2 NVIDIA Tesla K80 cards
- 1.8 (CPU) + 0.44 (GPU) PFLOP/s peak performance (#70)
- EDR InfiniBand



## JUHYDRA

- GPU prototyping machine
- 1 node with Intel Xeon E5 CPU ( $2 \times 8$  cores)
- NVIDIA  $2 \times$  Tesla K20,  $2 \times$  Tesla K40 cards
- No batch system

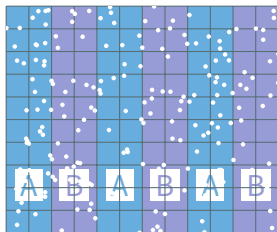


# JuSPIC

- Based on PSC by H. Ruhl
- Laser-plasma interaction
- 3D electromagnetic PiC code
- Finite-Difference  
Time-Domain scheme
- Cartesian geometry, arbitrary  
number of particle species
- Scales to full Blue Gene/Q  
system JUQUEEN

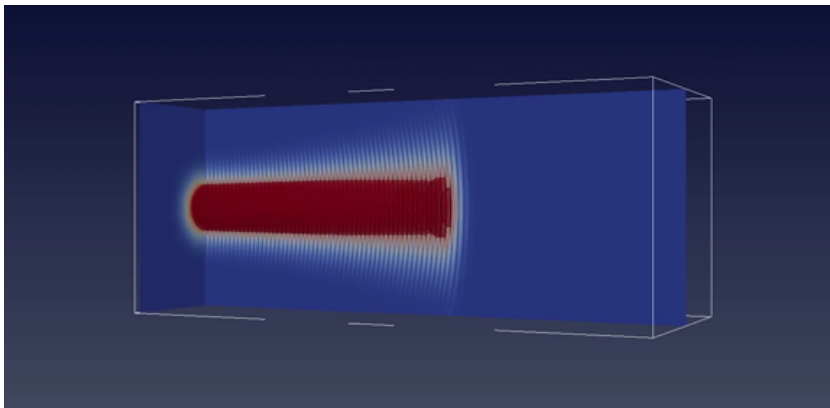


- Based on PSC by H. Ruhl
- Laser-plasma interaction
- 3D electromagnetic PiC code
- Finite-Difference Time-Domain scheme
- Cartesian geometry, arbitrary number of particle species
- Scales to full Blue Gene/Q system JUQUEEN
- Modern Fortran, Open Source
- Distributed with MPI in tiles
- CPU-parallelized with OpenMP



# Sample Simulation

*Visualizing different quantities*

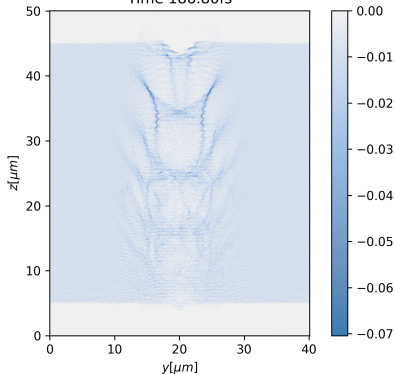


# Sample Simulation

*Visualizing different quantities*

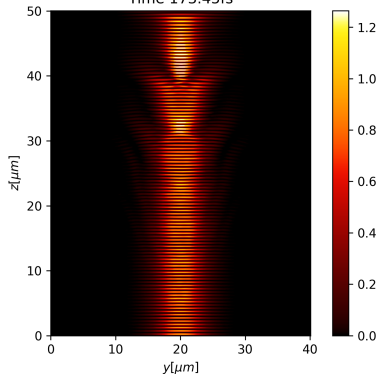
**Electron number density**

Time 186.80fs

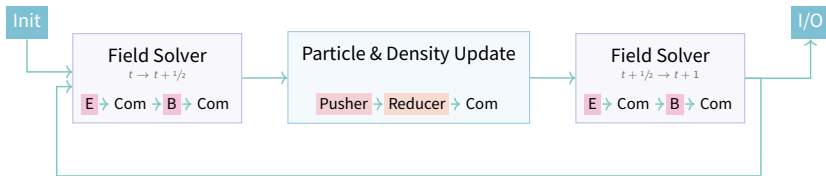


**E-Field  $E_y^2$**

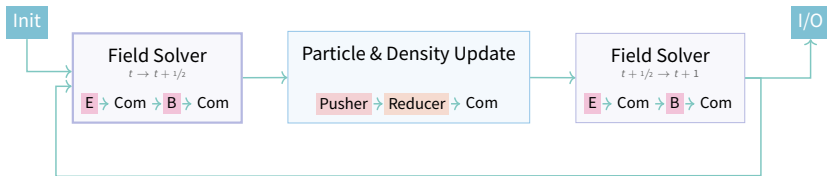
Time 173.45fs



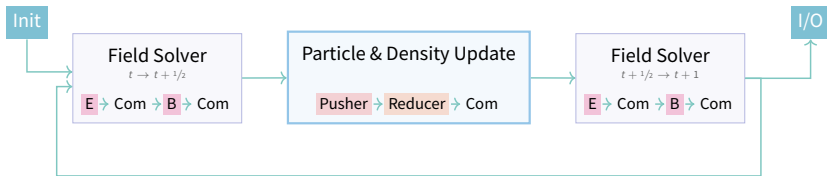
# Steps of Algorithm



# Steps of Algorithm

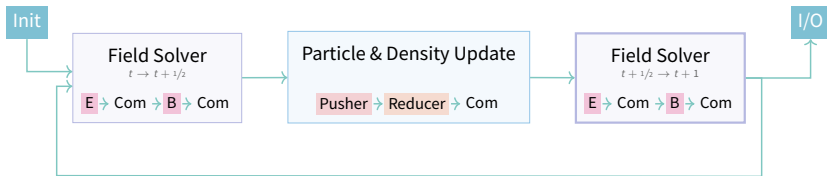


# Steps of Algorithm

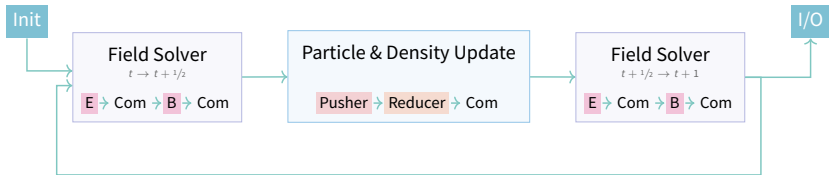




# Steps of Algorithm



# Steps of Algorithm



- **E, B** Already on GPU with OpenACC (small kernels)
- **Pusher** Focus of this paper
- **Reducer** Future step

# Acceleration for GPUs

# Acceleration for GPUs

## OpenACC

- Field solvers use OpenACC (simple code)

```
!$acc kernels loop collapse(3) present(e,b,ji)
do i3=i3mn-1,i3mx+1
  do i2=i2mn-1,i2mx+1
    do i1=i1mn-1,i1mx+1
      e(i1,i2,i3)%X=e(i1,i2,i3)%X
    ! etc
```

- Data movement with OpenACC (incl. resident parts)
- But Pusher no easy feat

- At start of porting: **Pusher** kernel too complicated for parsing by compiler
  - Large routine (many registers)
  - Operations on whole fields (it's Fortran after all)
  - Structured data types (with alloctables)
- Long investigation to get runnable code
- Good performance complicated
- Reported in other publication (beyond scope here, [appendix](#))



- At start of porting: **Pusher** kernel too complicated for parsing by compiler
    - Large routine (many registers)
    - Operations on whole fields (it's Fortran after all)
    - Structured data types (with alloctables)
  - Long investigation to get runnable code
  - Good performance complicated
  - Reported in other publication (beyond scope here, [appendix](#))
- Use CUDA Fortran



# Acceleration for GPUs

## CUDA Fortran

# Introduction to CUDA Fortran

*It's like CUDA C/C++,... but for Fortran*

- Available in PGI Fortran compiler
- Adds CUDA extensions to Fortran

# Introduction to CUDA Fortran

*It's like CUDA C/C++,... but for Fortran*

- Available in PGI Fortran compiler
- Adds CUDA extensions to Fortran
- Examples (from JuSPIC):

# Introduction to CUDA Fortran

*It's like CUDA C/C++,... but for Fortran*

- Available in PGI Fortran compiler
- Adds CUDA extensions to Fortran
- Examples (from JuSPIC):
  - Define device function along-side host function

```
type(particle_type), dimension(slice(1)%n) ::  
  ↪ list_of_particles, list_of_particles_d  
attributes(device) :: list_of_particles_d
```

# Introduction to CUDA Fortran

*It's like CUDA C/C++,... but for Fortran*

- Available in PGI Fortran compiler
- Adds CUDA extensions to Fortran
- Examples (from JuSPIC):
  - Define device function along-side host function

```
type(particle_type), dimension(slice(1)%n) ::  
  ↪ list_of_particles, list_of_particles_d  
attributes(device) :: list_of_particles_d
```

- Copy to device

```
list_of_particles_d = list_of_particles
```

# Introduction to CUDA Fortran

*It's like CUDA C/C++,... but for Fortran*

- Available in PGI Fortran compiler
- Adds CUDA extensions to Fortran
- Examples (from JuSPIC):
  - Define device function along-side host function

```
type(particle_type), dimension(slice(1)%n) ::  
  ↪ list_of_particles, list_of_particles_d  
attributes(device) :: list_of_particles_d
```

- Copy to device

```
list_of_particles_d = list_of_particles
```

- Define kernel

```
attributes(global) subroutine gpupusher(list_of_particles, ...)
```

# Introduction to CUDA Fortran

*It's like CUDA C/C++,... but for Fortran*

- Available in PGI Fortran compiler
- Adds CUDA extensions to Fortran
- Examples (from JuSPIC):

- Define device function along-side host function

```
type(particle_type), dimension(slice(1)%n) ::  
  ↪ list_of_particles, list_of_particles_d  
attributes(device) :: list_of_particles_d
```

- Copy to device

```
list_of_particles_d = list_of_particles
```

- Define kernel

```
attributes(global) subroutine gpupusher(list_of_particles, ...)
```

- Call kernel

```
call gpupusher<<<dim3(nBlocks, 1, 1), dim3(nThreads, 1,  
  ↪ 1)>>>(list_of_particles_d, ...)
```

# CUDA Fortran Portability

*Not as portable as OpenACC, but it's alright*

- CUDA Fortran: more powerful approach
- Portability suffers...



# CUDA Fortran Portability

*Not as portable as OpenACC, but it's alright*

- CUDA Fortran: more powerful approach
- Portability suffers...
- ... but can be mitigated!
  - 1 Use OpenACC as much as possible, e.g. for data movements  
*OpenACC mixes well together with CUDA Fortran*  
*!\$acc enter data copyin(list\_of\_particles, ...)*

# CUDA Fortran Portability

*Not as portable as OpenACC, but it's alright*

- CUDA Fortran: more powerful approach
  - Portability suffers...
  - ... but can be mitigated!
- 1 Use OpenACC as much as possible, e.g. for data movements

*OpenACC mixes well together with CUDA Fortran*

```
!$acc enter data copyin(list_of_particles, ...)
```

- 2 Use pre-processor directives for rest

```
#ifdef _CUDA
```

```
    i = blockDim%x * (blockIdx%x - 1) + threadIdx%x
```

```
#else
```

```
    do i = lbound(a, 1), ubound(a, 1)
```

```
#endif
```

# Acceleration for GPUs

## Data Layout Analysis

# Strategies for Data Layout

*Because data is not solely data*

- Benchmark different data layouts and transfer strategies
- Sub-parts of Pusher:

$\Sigma$  Everything  
**Allocate** Allocate host-side data structures  
**LL2F** Convert linked-list data structure to field  
**H2D** Copy data from host to device

**Kernel** Run kernel  
**D2H** Copy data from device to host  
**Other** Left-over time (synchronization, etc.)  
**F2LL** Copy flat field back to linked list

- Benchmarking on **JURON**

**Initial** All particles stored in single field, one particle after another; data copied to/from GPU with Fortran (*baseline*)

**Initial** All particles stored in single field, one particle after another; data copied to/from GPU with Fortran (*baseline*)

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |

**Initial** All particles stored in single field, one particle after another; data copied to/from GPU with Fortran (*baseline*)

**Exp 1** As *Initial*, but data copied with OpenACC *copy* directives

| <i>in</i> $\mu$ s | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-------------------|----------|----------|------|-----|--------|-----|--------|------|
| Initial           | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1             | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |

# Data Layout Experiments

## Description of experiments

**Initial** All particles stored in single field, one particle after another; data copied to/from GPU with Fortran (*baseline*)

**Exp 1** As *Initial*, but data copied with OpenACC *copy* directives

**Exp 2** As *Exp 1*, but data copied from pinned host memory

| <i>in</i> $\mu$ s | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-------------------|----------|----------|------|-----|--------|-----|--------|------|
| Initial           | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1             | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2             | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |



**Initial** All particles stored in single field, one particle after another; data copied to/from GPU with Fortran (*baseline*)

**Exp 1** As *Initial*, but data copied with OpenACC *copy* directives

**Exp 2** As *Exp 1*, but data copied from pinned host memory

**SoA** Data copied with Fortran, but instead of one field with all particle data, one field for each spatial and momentum component for particles



| <i>in</i> $\mu$ s | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-------------------|----------|----------|------|-----|--------|-----|--------|------|
| Initial           | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1             | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2             | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA               | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D       | Kernel    | D2H       | Others | F2LL |
|-----------------------------|----------|----------|------|-----------|-----------|-----------|--------|------|
| Initial                     | 8040     | –        | 567  | 82        | 84        | 62        | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80        | 82        | 91        | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79        | 83        | 72        | 108    | 7973 |
| <b>SoA</b>                  | 7811     | 1        | 844  | <b>66</b> | <b>77</b> | <b>53</b> | 376    | 6386 |

- **SoA**: fastest, looking (also) at raw GPU runtimes

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F       | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------------|-----|--------|-----|--------|------|
| Initial                     | 8040     | –        | 567        | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353        | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527        | 79  | 83     | 72  | 108    | 7973 |
| <b>SoA</b>                  | 7811     | 1        | <b>844</b> | 66  | 77     | 53  | 376    | 6386 |

- **SoA**: fastest, looking (also) at raw GPU runtimes – but slowest for change of data structures (six fields vs. one)

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others     | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|------------|------|
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350        | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380        | 9440 |
| <b>Exp 2</b>                | 9695     | 564      | 527  | 79  | 83     | 72  | <b>108</b> | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376        | 6386 |

- SoA: fastest, looking (also) at raw GPU runtimes – but slowest for change of data structures (six fields vs. one)
- **Exp 2**: least overhead; pinned memory allows for direct data access

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate   | LL2F | H2D | Kernel | D2H       | Others     | F2LL |
|-----------------------------|----------|------------|------|-----|--------|-----------|------------|------|
| Initial                     | 8040     | –          | 567  | 82  | 84     | 62        | 350        | 6885 |
| Exp 1                       | 10435    | –          | 353  | 80  | 82     | 91        | 380        | 9440 |
| <b>Exp 2</b>                | 9695     | <b>564</b> | 527  | 79  | 83     | <b>72</b> | <b>108</b> | 7973 |
| SoA                         | 7811     | 1          | 844  | 66  | 77     | 53        | 376        | 6386 |

- SoA: fastest, looking (also) at raw GPU runtimes – but slowest for change of data structures (six fields vs. one)
- **Exp 2**: least overhead; pinned memory allows for direct data access – but allocation overhead is not fully resolved

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL        |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|-------------|
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885        |
| <b>Exp 1</b>                | 10435    | –        | 353  | 80  | 82     | 91  | 380    | <b>9440</b> |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973        |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386        |

- SoA: fastest, looking (also) at raw GPU runtimes – but slowest for change of data structures (six fields vs. one)
- Exp 2: least overhead; pinned memory allows for direct data access – but allocation overhead is not fully resolved
- **Exp 1**: also ok for raw GPU times, but large F2LL overhead (*more on that later*)

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| <b>JURON</b>                |          |          |      |     |        |     |        |      |
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |
| <b>JUHYDRA</b>              |          |          |      |     |        |     |        |      |
| Initial                     | 4956     | –        | 908  | 267 | 229    | 208 | 736    | 2600 |
| Exp 1                       | 4687     | –        | 764  | 232 | 229    | 198 | 804    | 2455 |
| Exp 2                       | 5328     | 577      | 1027 | 224 | 230    | 192 | 23     | 2651 |
| SoA                         | 4880     | 1        | 786  | 204 | 208    | 173 | 827    | 2674 |



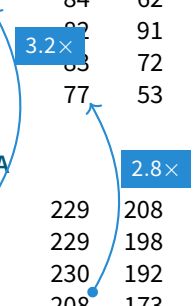
| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| <b>JURON</b>                |          |          |      |     |        |     |        |      |
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |
| <b>JUHYDRA</b>              |          |          |      |     |        |     |        |      |
| Initial                     | 4956     | –        | 908  | 267 | 229    | 208 | 736    | 2600 |
| Exp 1                       | 4687     | –        | 764  | 232 | 229    | 198 | 804    | 2455 |
| Exp 2                       | 5328     | 577      | 1027 | 224 | 230    | 192 | 23     | 2651 |
| SoA                         | 4880     | 1        | 786  | 204 | 208    | 173 | 827    | 2674 |

2.8×

# Data Layout Experiments

## Architecture Comparison

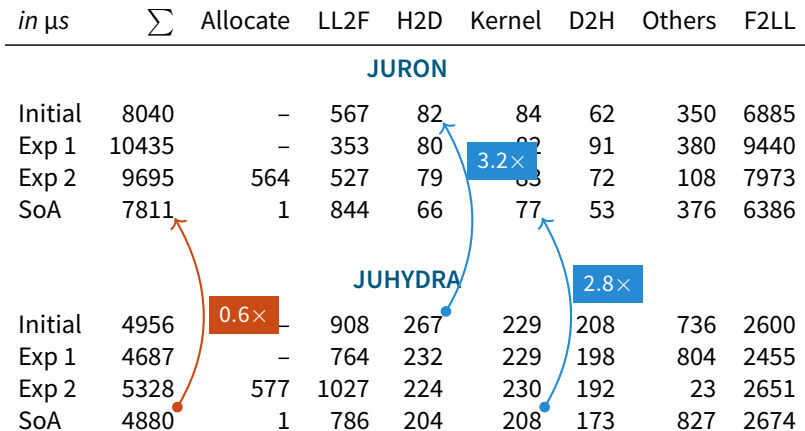
| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| <b>JURON</b>                |          |          |      |     |        |     |        |      |
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |
| <b>JUHYDRA</b>              |          |          |      |     |        |     |        |      |
| Initial                     | 4956     | –        | 908  | 267 | 229    | 208 | 736    | 2600 |
| Exp 1                       | 4687     | –        | 764  | 232 | 229    | 198 | 804    | 2455 |
| Exp 2                       | 5328     | 577      | 1027 | 224 | 230    | 192 | 23     | 2651 |
| SoA                         | 4880     | 1        | 786  | 204 | 208    | 173 | 827    | 2674 |



# Data Layout Experiments

## Architecture Comparison

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| <b>JURON</b>                |          |          |      |     |        |     |        |      |
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |
| <b>JUHYDRA</b>              |          |          |      |     |        |     |        |      |
| Initial                     | 4956     | 0.6×     | 908  | 267 | 229    | 208 | 736    | 2600 |
| Exp 1                       | 4687     | –        | 764  | 232 | 229    | 198 | 804    | 2455 |
| Exp 2                       | 5328     | 577      | 1027 | 224 | 230    | 192 | 23     | 2651 |
| SoA                         | 4880     | 1        | 786  | 204 | 208    | 173 | 827    | 2674 |



# Data Layout Experiments

## Architecture Comparison

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| <b>JURON</b>                |          |          |      |     |        |     |        |      |
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    | –        | 353  | 80  | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     | 564      | 527  | 79  | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     | 1        | 844  | 66  | 77     | 53  | 376    | 6386 |
| <b>JUHYDRA</b>              |          |          |      |     |        |     |        |      |
| Initial                     | 4956     | 0.6×     | 908  | 267 | 229    | 208 | 736    | 2600 |
| Exp 1                       | 4687     | –        | 764  | 232 | 229    | 198 | 804    | 2455 |
| Exp 2                       | 5328     | 577      | 1027 | 224 | 230    | 192 | 23     | 2651 |
| SoA                         | 4880     | 1        | 786  | 204 | 208    | 173 | 827    | 2674 |

# Data Layout Experiments

## Architecture Comparison

| <i>in <math>\mu</math>s</i> | $\Sigma$ | Allocate | LL2F | H2D | Kernel | D2H | Others | F2LL |
|-----------------------------|----------|----------|------|-----|--------|-----|--------|------|
| <b>JURON</b>                |          |          |      |     |        |     |        |      |
| Initial                     | 8040     | –        | 567  | 82  | 84     | 62  | 350    | 6885 |
| Exp 1                       | 10435    |          |      |     | 82     | 91  | 380    | 9440 |
| Exp 2                       | 9695     |          |      |     | 83     | 72  | 108    | 7973 |
| SoA                         | 7811     |          |      |     | 77     | 53  | 376    | 6386 |
| <b>JUHYDRA</b>              |          |          |      |     |        |     |        |      |
| Initial                     | 4956     | 0.6×     | 908  | 267 | 229    | 208 | 736    | 2600 |
| Exp 1                       | 4687     | –        | 764  | 232 | 229    | 198 | 804    | 2455 |
| Exp 2                       | 5328     | 577      | 1027 | 224 | 230    | 192 | 23     | 2651 |
| SoA                         | 4880     | 1        | 786  | 204 | 208    | 173 | 827    | 2674 |

Why!?

0.6×

2×

2.8×

0.3×

# Acceleration for GPUs

## Data Layout Conversion

- Parts of F2LL
  - Kill old linked list of particles<sup>1</sup>
  - Initialize new, empty linked list of particles
  - Loop through field(s) of particle information...
  - ... add each particle to linked list, update pointers

---

<sup>1</sup>Start with first particle, progress along links, remove each particle

# Conversion of Data Layouts

*Why is F2LL so slow?*

- Parts of F2LL
  - Kill old linked list of particles<sup>1</sup>
  - Initialize new, empty linked list of particles
  - Loop through field(s) of particle information...
  - ... add each particle to linked list, update pointers

- `add_one_to_list`

```
allocate(list%tail%next)
nullify(list%tail%next%next)
list%tail%next%particle = particle
list%tail => list%tail%next
```

---

<sup>1</sup>Start with first particle, progress along links, remove each particle



# Conversion of Data Layouts

*Why is F2LL so slow?*

- Parts of F2LL
  - Kill old linked list of particles<sup>1</sup>
  - Initialize new, empty linked list of particles
  - Loop through field(s) of particle information...
  - ... add each particle to linked list, update pointers

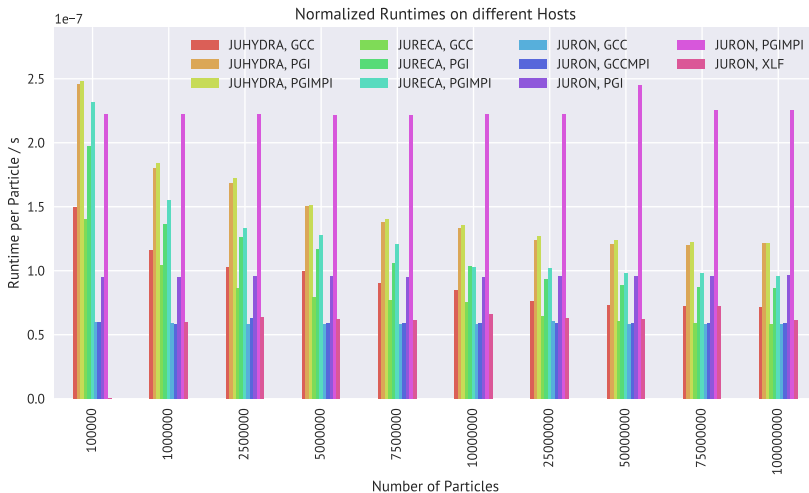
- `add_one_to_list`

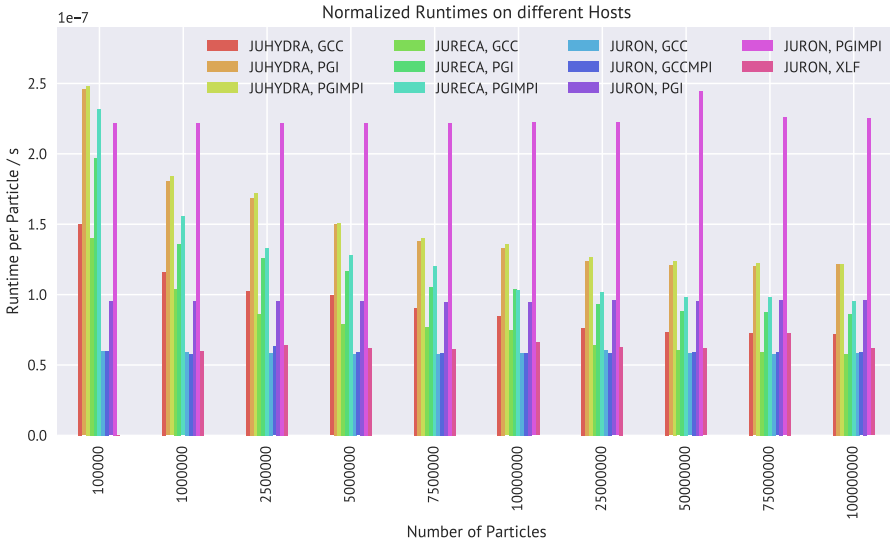
```
allocate(list%tail%next)
nullify(list%tail%next%next)
list%tail%next%particle = particle
list%tail => list%tail%next
```

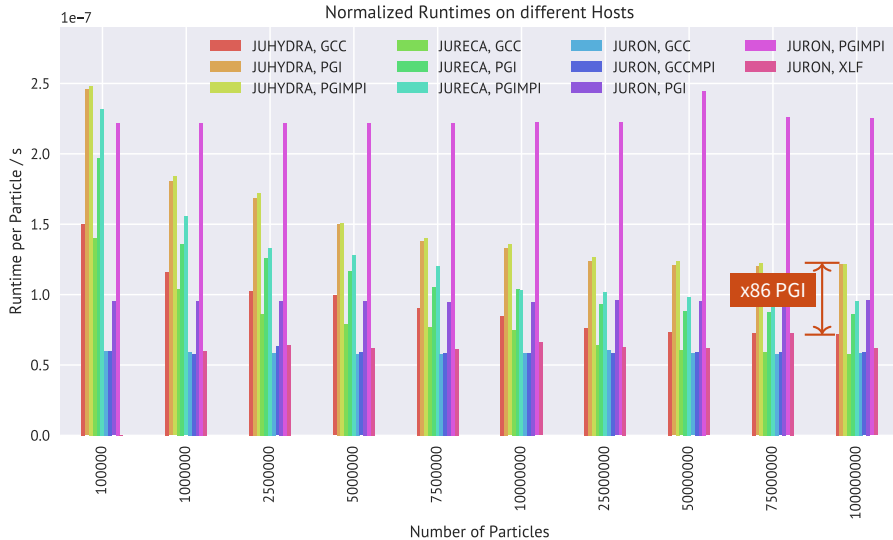
⇒ *Benchmark*

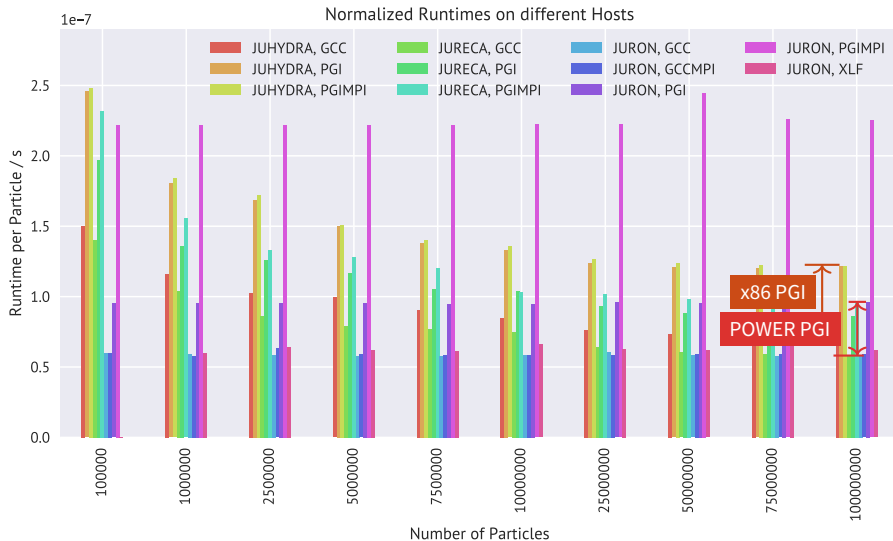
---

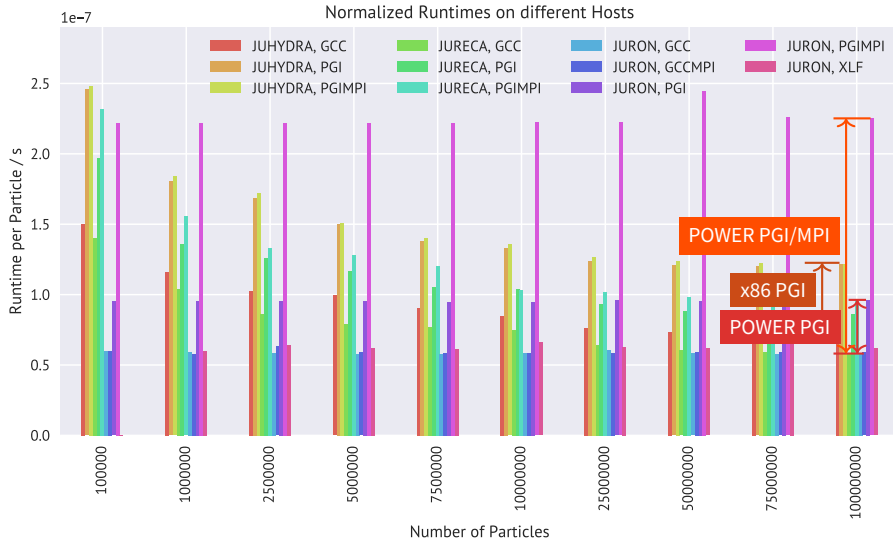
<sup>1</sup>Start with first particle, progress along links, remove each particle











# Compiler Investigation

*Is MPI Slow? And, by the way, which MPI!?*

- *PGIMPI*: MPI version shipped with PGI
- Not actively used in GPU version of JuSPIC, but in future

- *PGIMPI*: MPI version shipped with PGI
- Not actively used in GPU version of JuSPIC, but in future
- `add_one_to_list` benchmark does not use MPI at all!  
Replacing `pgfortran` by `mpifort` leads to performance decrease



- *PGIMPI*: MPI version shipped with PGI
  - Not actively used in GPU version of JuSPIC, but in future
  - `add_one_to_list` benchmark does not use MPI at all!  
Replacing `pgfortran` by `mpifort` leads to performance decrease
- **Benchmark** compilers – with PAPI [3] instrumentation

- *PGIMPI*: MPI version shipped with PGI
  - Not actively used in GPU version of JuSPIC, but in future
  - `add_one_to_list` benchmark does not use MPI at all!  
Replacing `pgfortran` by `mpifort` leads to performance decrease
- **Benchmark** compilers – with PAPI [3] instrumentation

| System          | JURON |        |     |        |         |     | JUHYDRA |        |
|-----------------|-------|--------|-----|--------|---------|-----|---------|--------|
| Compiler        | GCC   | GCCMPI | PGI | PGIMPI | PGIMPI* | XLF | PGI     | PGIMPI |
| Time pP/ns      | 36    | 37     | 46  | 154    | 48      | 41  | 32      | 32     |
| Instructions pP | 121   | 121    | 243 | 462    | 243     | 121 | 210     | 210    |

See [appendix](#) for some more counters

- MPI version shipped with PGI on POWER is slow, because it issues many instructions

- MPI version shipped with PGI on POWER is slow, because it issues many instructions
- Further study: Identical assembly code generated as MPI-less version...

- MPI version shipped with PGI on POWER is slow, because it issues many instructions
- Further study: Identical assembly code generated as MPI-less version...
- ... but includes call to `malloc()`!

- MPI version shipped with PGI on POWER is slow, because it issues many instructions
- Further study: Identical assembly code generated as MPI-less version...
- ... but includes call to `malloc()`!
- Different libraries linked for PGI and PGIMPI cases!

- MPI version shipped with PGI on POWER is slow, because it issues many instructions
- Further study: Identical assembly code generated as MPI-less version...
- ... but includes call to `malloc()`!
- Different libraries linked for PGI and PGIMPI cases!
- `LD_PRELOAD=/lib64/libc.so.6` solves problem!

- MPI version shipped with PGI on POWER is slow, because it issues many instructions
  - Further study: Identical assembly code generated as MPI-less version...
  - ... but includes call to `malloc()`!
  - Different libraries linked for PGI and PGIMPI cases!
  - `LD_PRELOAD=/lib64/libc.so.6` solves problem!
- ⇒ Slow MPI-aware `malloc()`?



- MPI version shipped with PGI on POWER is slow, because it issues many instructions
  - Further study: Identical assembly code generated as MPI-less version...
  - ... but includes call to `malloc()`!
  - Different libraries linked for PGI and PGIMPI cases!
  - `LD_PRELOAD=/lib64/libc.so.6` solves problem!
- ⇒ Slow MPI-aware `malloc()`?
- Mitigation
    - **Bug reported**
    - For now: consider as *anomalous* overhead

# Performance Modelling

- **Goal:** Compare different GPU architectures; understand behavior of JuSPIC
- Model based on **information exchanged** of GPU kernel
  - Amount of exchanged information for given number of particles
  - Time for exchange

# Effective Bandwidth

## *Defining the model*

- **Goal:** Compare different GPU architectures; understand behavior of JuSPIC
- Model based on **information exchanged** of GPU kernel
  - Amount of exchanged information for given number of particles
  - Time for exchange

$$t(N_{\text{part}}) = \alpha + I(N_{\text{part}})/\beta ,$$

- **Goal:** Compare different GPU architectures; understand behavior of JuSPIC
- Model based on **information exchanged** of GPU kernel
  - Amount of exchanged information for given number of particles
  - Time for exchange

$$t(N_{\text{part}}) = \alpha + I(N_{\text{part}})/\beta ,$$

$N_{\text{part}}$  Number of particles processed

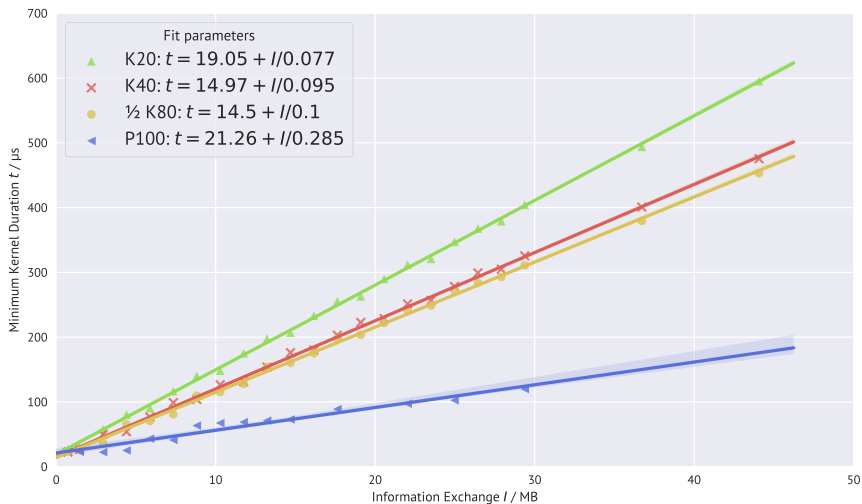
$I$  Information exchanged (572 B (read) + 40 B (write))

$t$  Kernel runtime

$\alpha, \beta$  Fit parameters;  $\beta$ : *effective bandwidth*

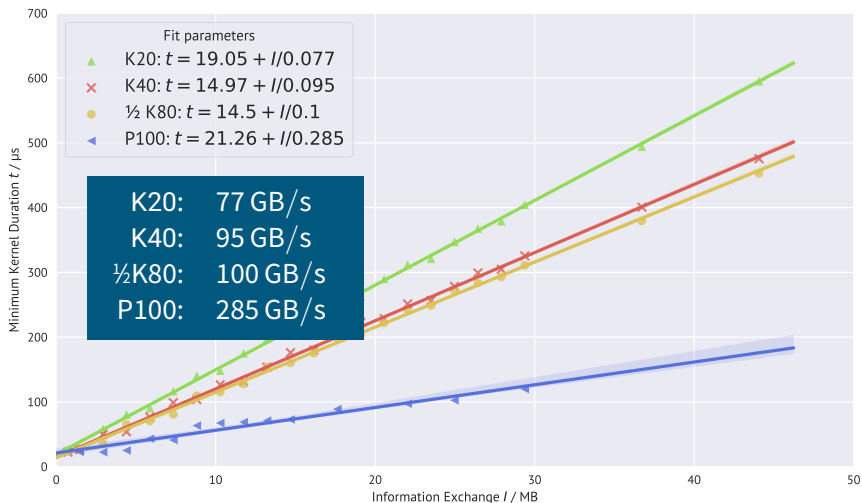
# Effective Bandwidth

## Measurements



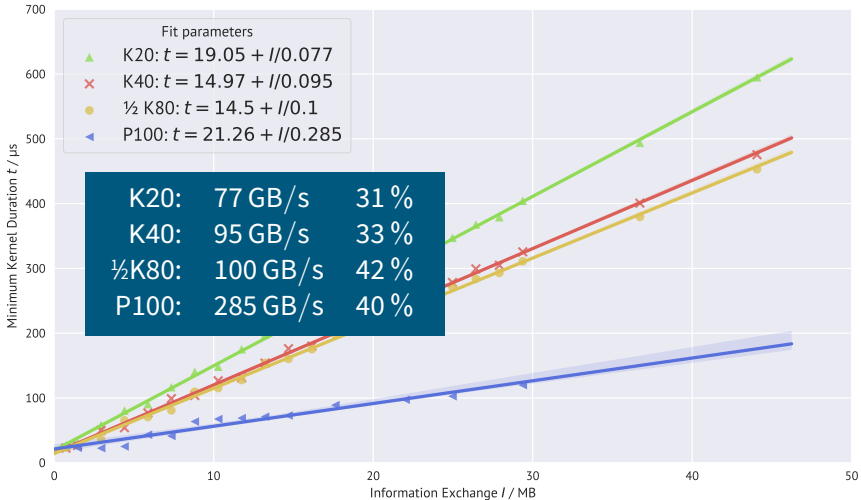
# Effective Bandwidth

## Measurements



# Effective Bandwidth

## Measurements





- Another free parameter: **GPU clock rates**
  - Varies significantly across GPU architecture generations and models
- Incorporate clock into performance model

- Another free parameter: **GPU clock rates**
  - Varies significantly across GPU architecture generations and models
- Incorporate clock into performance model

$$\beta(\mathcal{C}) = \gamma + \delta \mathcal{C}$$

# Clock Dependency

## *Defining the relation*

- Another free parameter: **GPU clock rates**
  - Varies significantly across GPU architecture generations and models
- Incorporate clock into performance model

$$\beta(\mathcal{C}) = \gamma + \delta \mathcal{C}$$

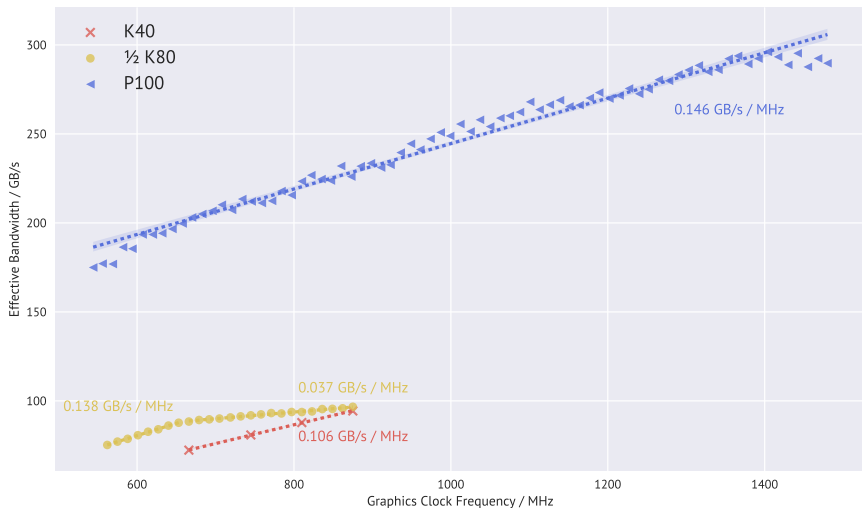
$\mathcal{C}$  GPU clock rate

$\beta$  Effective bandwidth (from before)

$\gamma, \delta$  Fit parameters

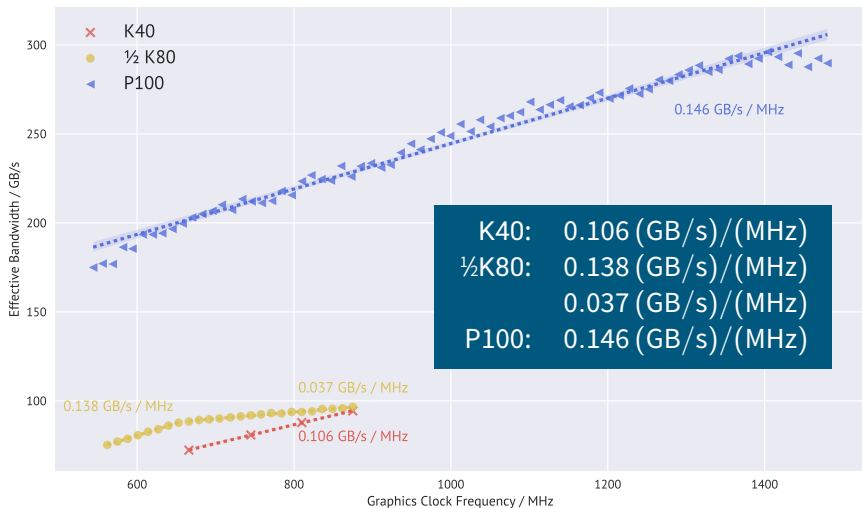
# Clock Dependency

## Measurements



# Clock Dependency

## Measurements



## Summary

- Enabled **JuSPIC** for **GPU** with OpenACC & CUDA Fortran
- Particle data layout: **SoA fastest**
- **Slow memory allocation** with PGI+MPI on POWER → bug filed
- Performance model: **Information exchange** (P100: 285 GB/s)
- Studied model with **different clock rates** – P100 most efficient scaling

## Future

- Port also Reducer to GPU
- Enable MPI again
- Alternatives to linked list

## Summary

- Enabled **JuSPIC** for **GPU** with OpenACC & CUDA Fortran
- Particle data layout: **SoA fastest**
- **Slow memory allocation** with PGI+MPI on POWER → bug filed
- Performance model: **Information exchange** (P100: 285 GB/s)
- Studied model with **different clock rates** – P100 most efficient scaling

## Future

- Port also Reducer to GPU
- Enable MPI again
- Alternatives to linked list

*Thank you  
for your attention!*  
*a.herten@fz-juelich.de*

## Appendix

Acknowledgements

Related Work

OpenACC Performance Progression

Linked List: Remove on JURON

Selected Performance Counters on JURON

References

Glossary



- The work was done in context of two groups:  
**POWER Acceleration and Design Centre** A collaboration of IBM, NVIDIA, and Forschungszentrum Jülich  
**NVIDIA Application Lab** A collaboration of NVIDIA and Forschungszentrum Jülich
- Many thanks to Jiri Kraus from NVIDIA, who helped tremendously along the way
- JURON, a prototype system for the Human Brain Project, received co-funding from the European Union (Grant Agreement No. 604102)

- Selection of other GPU PiC codes
  - PSC The code JuSPIC is based on has been reimplemented in C and ported to GPU [4]
  - PIConGPU PiC code specifically developed for GPUs [5]
- Minsky porting experiences
  - “Addressing Materials Science Challenges Using GPU-accelerated POWER8 Nodes” [6]
  - “A Performance Model for GPU-Accelerated FDTD Applications” [7]
- ... more in paper!

# OpenACC Performance Progression

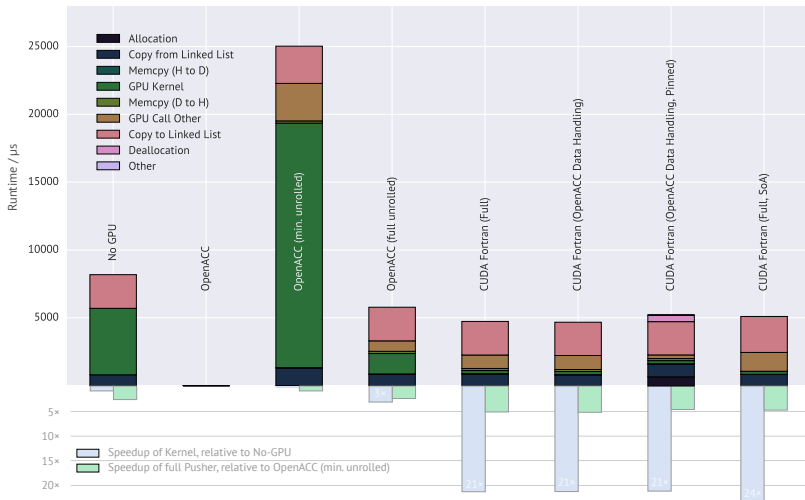
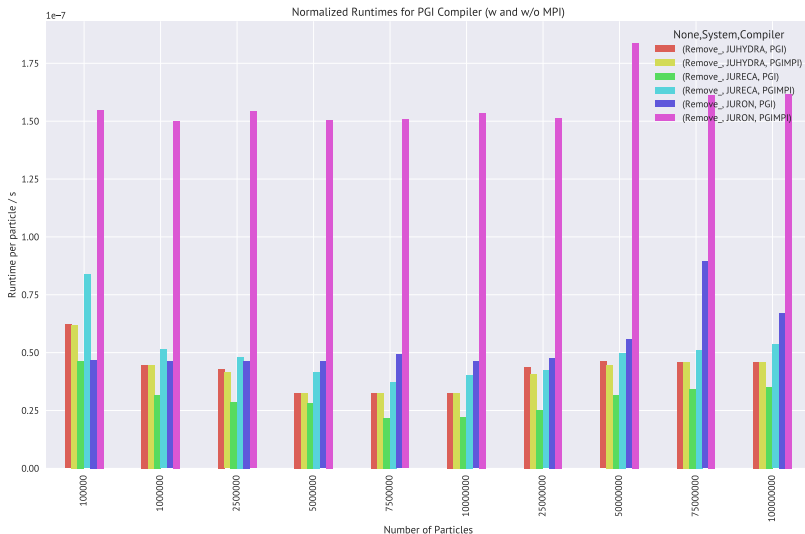


Figure: See GTC poster for details [8].

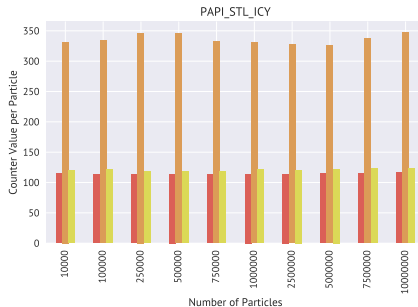
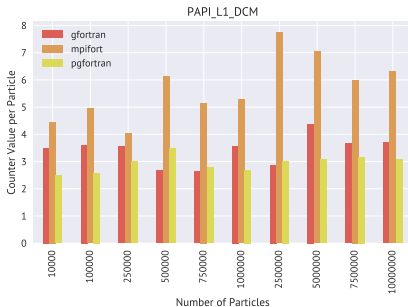
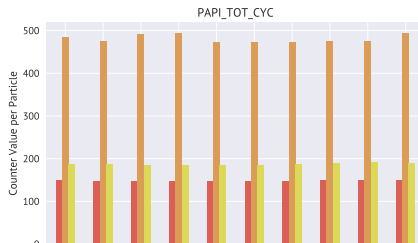
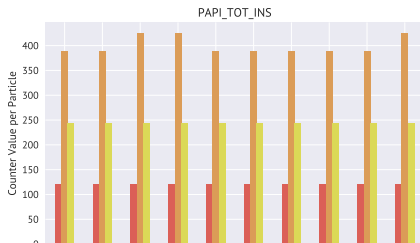
# Linked List: Time for Remove on JURON

*For different compilers*



# Selected Performance Counters on JURON

*For different compilers*



- [1] Forschungszentrum Jülich. *Hightech made in 1960: A view into the control room of DIDO*. URL: [http://historie.fz-juelich.de/60jahre/DE/Geschichte/1956-1960/Dekade/\\_node.html](http://historie.fz-juelich.de/60jahre/DE/Geschichte/1956-1960/Dekade/_node.html) (page 3).
- [2] Forschungszentrum Jülich. *Forschungszentrum Bird's Eye*. (Page 3).

- [3] Phil Mucci and The ICL Team. *PAPI, the Performance Application Programming Interface*. URL: <http://icl.utk.edu/papi/> (visited on 04/30/2017) (pages 63–66, 100).
- [4] K. Germaschewski et al. “The Plasma Simulation Code: A modern particle-in-cell code with load-balancing and GPU support”. In: *ArXiv e-prints* (Oct. 2013). arXiv: 1310.7866 [physics.plasm-ph] (page 90).
- [5] M. Bussmann et al. “Radiative signature of the relativistic Kelvin-Helmholtz Instability”. In: *2013 SC - International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*. Nov. 2013, pp. 1–12. DOI: 10.1145/2503210.2504564 (page 90).

- [6] Paul F. Baumeister et al. “Addressing Materials Science Challenges Using GPU-accelerated POWER8 Nodes”. In: *Euro-Par 2016: Parallel Processing: 22nd International Conference on Parallel and Distributed Computing, Grenoble, France, August 24-26, 2016, Proceedings*. Ed. by Pierre-François Dutot and Denis Trystram. Cham: Springer International Publishing, 2016, pp. 77–89. ISBN: 978-3-319-43659-3. DOI: 10.1007/978-3-319-43659-3\_6. URL: [http://dx.doi.org/10.1007/978-3-319-43659-3\\_6](http://dx.doi.org/10.1007/978-3-319-43659-3_6) (page 90).



- [7] P. F. Baumeister et al. “A Performance Model for GPU-Accelerated FDTD Applications”. In: *2015 IEEE 22nd International Conference on High Performance Computing (HiPC)*. Dec. 2015, pp. 185–193. DOI: 10.1109/HiPC.2015.24 (page 90).
- [8] Andreas Herten, Dirk Pleiter, and Dirk Brömmel. *Accelerating Plasma Physics with GPUs (Poster)*. Tech. rep. GPU Technology Conference, 2017 (page 91).
- [9] Philip J. Mucci et al. “PAPI: A Portable Interface to Hardware Performance Counters”. In: *In Proceedings of the Department of Defense HPCMP Users Group Conference*. 1999, pp. 7–10 (page 100).

**CUDA** Computing platform for **GPUs** from NVIDIA. Provides, among others, CUDA C/C++. [2](#), [22](#), [23](#), [24](#), [26](#), [27](#), [28](#), [29](#), [30](#), [31](#), [32](#), [33](#), [34](#), [86](#), [87](#)

**FZJ** Forschungszentrum Jülich, a research center in the west of Germany. [3](#), [98](#)

**JSC** Jülich Supercomputing Centre operates a number of large and small supercomputers and connected infrastructure at **FZJ**. [3](#)

**JuSPIC** Jülich Scalable Particle-in-Cell Code. [2](#), [9](#), [10](#), [11](#), [26](#), [27](#), [28](#), [29](#), [30](#), [31](#), [63](#), [64](#), [65](#), [66](#), [75](#), [76](#), [77](#), [86](#), [87](#), [90](#)

**MPI** The Message Passing Interface, a communication message-passing application programmer interface.  
63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 86, 87

**NVIDIA** US technology company creating GPUs. 3, 4, 5, 6, 7, 8, 89, 98

**NVLink** NVIDIA's communication protocol connecting CPU ↔ GPU and GPU ↔ GPU with 80 GB/s. PCI-Express: 16 GB/s. 4, 5, 6, 7, 8, 98

**OpenACC** Directive-based programming, primarily for many-core machines. 2, 14, 15, 16, 17, 18, 21, 22, 23, 24, 32, 33, 34, 37, 38, 39, 40, 41, 86, 87, 88, 91

- P100** A large GPU with the Pascal architecture from NVIDIA. It employs NVLink as its interconnect and has fast HBM2 memory. 2, 4, 5, 6, 7, 8, 79, 80, 85, 86, 87
- PAPI** The Performance API, a interface for accessing performance counters, also with aliased names cross-platform [3, 9]. 63, 64, 65, 66
- Pascal** The latest available GPU architecture from NVIDIA. 98
- PGI** Formerly *The Portland Group, Inc.*; since 2013 part of NVIDIA. 2, 26, 27, 28, 29, 30, 31, 67, 68, 69, 70, 71, 72, 73, 86, 87
- PiC** Particle in Cell; a method applied in a group of (plasma) physics simulations to solve partial differential equations. 2, 10, 11, 90

**POWER** Series of microprocessors from IBM. 2, 3, 67, 68, 69, 70, 71, 72, 73, 86, 87, 89

**Tesla** The GPU product line for general purpose computing computing of NVIDIA. 4, 5, 6, 7, 8