



University of Stuttgart
Germany



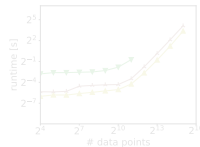
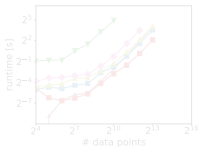
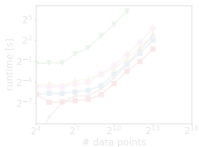
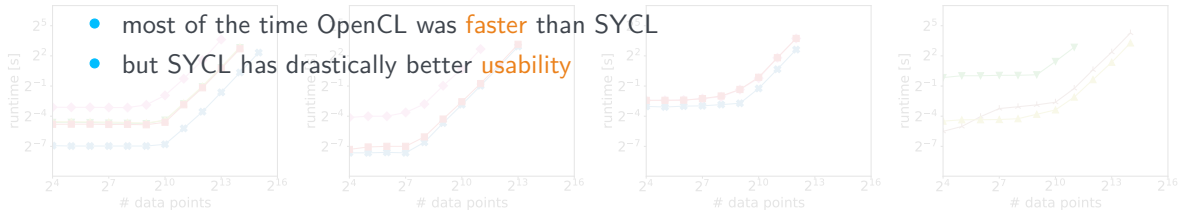
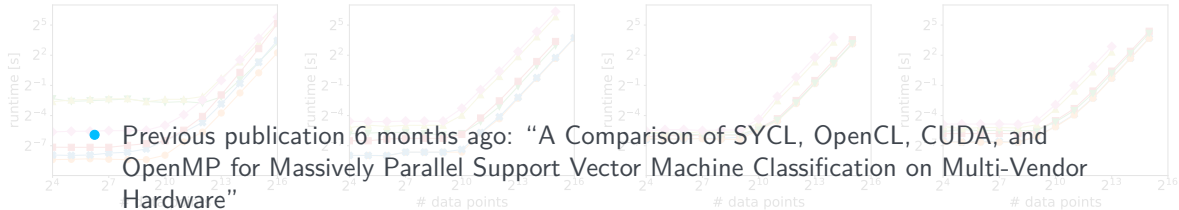
Marcel
Breyer

*oneAPI DevSummit
for HPC and AI*

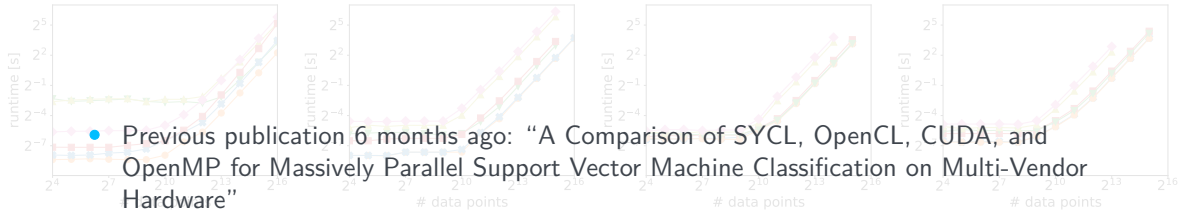
6-7 December 2022

**Performance Evolution
of Different SYCL
Implementations on
the Basis of PLSSVM**

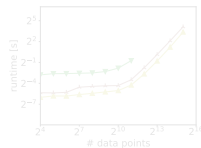
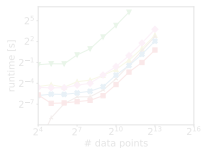
Motivation



Motivation



Could SYCL close the performance gap to OpenCL in our use case?

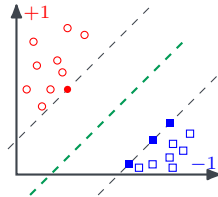


**What
to know
about
PLSSVM**

1

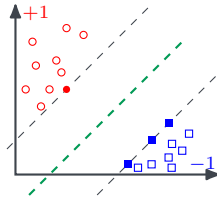
Support Vector Machines (SVMs) and their problems

- SVMs as supervised machine learning technique
- originally meant for binary **classification**



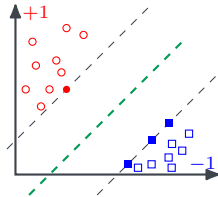
Support Vector Machines (SVMs) and their problems

- SVMs as supervised machine learning technique
- originally meant for binary **classification**
- SVMs have to solve a convex quadratic problem
 - state-of-the-art: Sequential Minimal Optimization (SMO) (proposed by Platt in 1998)
 - **inherently sequential algorithm**



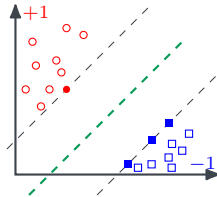
Support Vector Machines (SVMs) and their problems

- SVMs as supervised machine learning technique
- originally meant for binary **classification**
- SVMs have to solve a convex quadratic problem
 - state-of-the-art: Sequential Minimal Optimization (SMO) (proposed by Platt in 1998)
 - **inherently sequential algorithm**
- many SVM implementations modify SMO to exploit some parallelism
 - still not well suited for modern, highly parallel hardware



Support Vector Machines (SVMs) and their problems

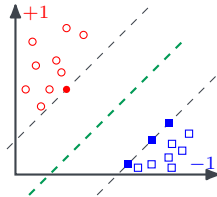
- SVMs as supervised machine learning technique
- originally meant for binary **classification**
- SVMs have to solve a convex quadratic problem
 - state-of-the-art: Sequential Minimal Optimization (SMO) (proposed by Platt in 1998)
 - **inherently sequential algorithm**
- many SVM implementations modify SMO to exploit some parallelism
 - still not well suited for modern, highly parallel hardware



→ *Least Squares Support Vector Machine (LS-SVM)*

(proposed by Suykens and Vandewalle in 1999)

Support Vector Machines (SVMs) and their problems



- SVMs as supervised machine learning technique
- originally meant for binary **classification**
- SVMs have to solve a convex quadratic problem
 - state-of-the-art: Sequential Minimal Optimization (SMO) (proposed by Platt in 1998)
 - **inherently sequential algorithm**
- many SVM implementations modify SMO to exploit some parallelism
 - still not well suited for modern, highly parallel hardware

→ *Least Squares Support Vector Machine (LS-SVM)*

(proposed by Suykens and Vandewalle in 1999)

- reformulation of standard SVM to **solving a system of linear equations**
- massively parallel algorithms known

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} \mathbf{Q} & \vec{\mathbf{1}}_n \\ \vec{\mathbf{1}}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \boldsymbol{\alpha} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}$$

where \mathbf{Q} is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is symmetric positive-definite

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is **symmetric positive-definite**

→ Conjugate Gradient algorithm: (variant of Shewchuk et al.)

```
1:  $i \leftarrow 0$ 
2:  $r \leftarrow b - Ax$ 
3:  $d \leftarrow r$ 
4:  $\delta_{new} \leftarrow r^T r$ 
5:  $\delta_0 \leftarrow \delta_{new}$ 
6: while  $i < i_{max}$  and  $\delta_{new} > \epsilon^2 \delta_0$  do
7:    $q \leftarrow Ad$ 
8:    $\alpha \leftarrow \frac{\delta_{new}}{d^T q}$ 
9:    $x \leftarrow x + \alpha d$ 
10:  if  $i$  is divisible by 50 then
11:     $r \leftarrow b - Ax$ 
12:  else
13:     $r \leftarrow r - \alpha q$ 
14:  end if
15:   $\delta_{old} \leftarrow \delta_{new}$ 
16:   $\delta_{new} \leftarrow r^T r$ 
17:   $\beta \leftarrow \frac{\delta_{new}}{\delta_{old}}$ 
18:   $d \leftarrow r + \beta d$ 
19:   $i \leftarrow i + 1$ 
20: end while
```

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is **symmetric positive-definite**

→ Conjugate Gradient algorithm: (variant of Shewchuk et al.)

- **Setup or constant operations** → host

```
1:  $i \leftarrow 0$ 
2:  $r \leftarrow b - Ax$ 
3:  $d \leftarrow r$ 
4:  $\delta_{new} \leftarrow r^T r$ 
5:  $\delta_0 \leftarrow \delta_{new}$ 
6: while  $i < i_{max}$  and  $\delta_{new} > \epsilon^2 \delta_0$  do
7:    $q \leftarrow Ad$ 
8:    $\alpha \leftarrow \frac{\delta_{new}}{d^T q}$ 
9:    $x \leftarrow x + \alpha d$ 
10:  if  $i$  is divisible by 50 then
11:     $r \leftarrow b - Ax$ 
12:  else
13:     $r \leftarrow r - \alpha q$ 
14:  end if
15:   $\delta_{old} \leftarrow \delta_{new}$ 
16:   $\delta_{new} \leftarrow r^T r$ 
17:   $\beta \leftarrow \frac{\delta_{new}}{\delta_{old}}$ 
18:   $d \leftarrow r + \beta d$ 
19:   $i \leftarrow i + 1$ 
20: end while
```

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is symmetric positive-definite

→ Conjugate Gradient algorithm: (variant of Shewchuk et al.)

- Setup or constant operations → host
- BLAS Level 1 → host

```
1:  $i \leftarrow 0$ 
2:  $r \leftarrow b - Ax$ 
3:  $d \leftarrow r$ 
4:  $\delta_{new} \leftarrow r^T r$ 
5:  $\delta_0 \leftarrow \delta_{new}$ 
6: while  $i < i_{max}$  and  $\delta_{new} > \epsilon^2 \delta_0$  do
7:    $q \leftarrow Ad$ 
8:    $\alpha \leftarrow \frac{\delta_{new}}{d^T q}$ 
9:    $x \leftarrow x + \alpha d$ 
10:  if  $i$  is divisible by 50 then
11:     $r \leftarrow b - Ax$ 
12:  else
13:     $r \leftarrow r - \alpha q$ 
14:  end if
15:   $\delta_{old} \leftarrow \delta_{new}$ 
16:   $\delta_{new} \leftarrow r^T r$ 
17:   $\beta \leftarrow \frac{\delta_{new}}{\delta_{old}}$ 
18:   $d \leftarrow r + \beta d$ 
19:   $i \leftarrow i + 1$ 
20: end while
```

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is symmetric positive-definite

→ Conjugate Gradient algorithm: (variant of Shewchuk et al.)

- Setup or constant operations → host
- BLAS Level 1 → host
- BLAS Level 2 → device

```
1:  $i \leftarrow 0$ 
2:  $r \leftarrow b - Ax$ 
3:  $d \leftarrow r$ 
4:  $\delta_{new} \leftarrow r^T r$ 
5:  $\delta_0 \leftarrow \delta_{new}$ 
6: while  $i < i_{max}$  and  $\delta_{new} > \epsilon^2 \delta_0$  do
7:    $q \leftarrow Ad$ 
8:    $\alpha \leftarrow \frac{\delta_{new}}{d^T q}$ 
9:    $x \leftarrow x + \alpha d$ 
10:  if  $i$  is divisible by 50 then
11:     $r \leftarrow b - Ax$ 
12:  else
13:     $r \leftarrow r - \alpha q$ 
14:  end if
15:   $\delta_{old} \leftarrow \delta_{new}$ 
16:   $\delta_{new} \leftarrow r^T r$ 
17:   $\beta \leftarrow \frac{\delta_{new}}{\delta_{old}}$ 
18:   $d \leftarrow r + \beta d$ 
19:   $i \leftarrow i + 1$ 
20: end while
```

We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is **symmetric positive-definite**

→ Conjugate Gradient algorithm: (variant of Shewchuk et al.)

- **Setup or constant operations** → host
- **BLAS Level 1** → host
- **BLAS Level 2** → device

→ $Q \in \mathbb{R}^{\text{num_data_points} \times \text{num_data_points}}$

```
1:  $i \leftarrow 0$ 
2:  $r \leftarrow b - Ax$ 
3:  $d \leftarrow r$ 
4:  $\delta_{new} \leftarrow r^T r$ 
5:  $\delta_0 \leftarrow \delta_{new}$ 
6: while  $i < i_{max}$  and  $\delta_{new} > \epsilon^2 \delta_0$  do
7:    $q \leftarrow Ad$ 
8:    $\alpha \leftarrow \frac{\delta_{new}}{d^T q}$ 
9:    $x \leftarrow x + \alpha d$ 
10:  if  $i$  is divisible by 50 then
11:     $r \leftarrow b - Ax$ 
12:  else
13:     $r \leftarrow r - \alpha q$ 
14:  end if
15:   $\delta_{old} \leftarrow \delta_{new}$ 
16:   $\delta_{new} \leftarrow r^T r$ 
17:   $\beta \leftarrow \frac{\delta_{new}}{\delta_{old}}$ 
18:   $d \leftarrow r + \beta d$ 
19:   $i \leftarrow i + 1$ 
20: end while
```


We parallelized the most complex operations in the CG algorithm

LS-SVMs solve the system of linear equations:

$$\begin{bmatrix} Q & \vec{1}_n \\ \vec{1}_n^T & 0 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ b \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

where Q is the kernel matrix according to

$$Q_{ij} = k(\vec{x}_i, \vec{x}_j) + \frac{1}{C} \cdot \delta_{ij} \quad \left(\text{with } \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

→ Q is **symmetric positive-definite**

→ Conjugate Gradient algorithm: (variant of Shewchuk et al.)

- **Setup or constant operations** → host
- **BLAS Level 1** → host
- **BLAS Level 2** → device

→ $Q \in \mathbb{R}^{\text{num_data_points} \times \text{num_data_points}}$

→ **implicitly** calculate Q in each iteration

```
1:  $i \leftarrow 0$ 
2:  $r \leftarrow b - Ax$ 
3:  $d \leftarrow r$ 
4:  $\delta_{new} \leftarrow r^T r$ 
5:  $\delta_0 \leftarrow \delta_{new}$ 
6: while  $i < i_{max}$  and  $\delta_{new} > \epsilon^2 \delta_0$  do
7:    $q \leftarrow Ad$ 
8:    $\alpha \leftarrow \frac{\delta_{new}}{d^T q}$ 
9:    $x \leftarrow x + \alpha d$ 
10:  if  $i$  is divisible by 50 then
11:     $r \leftarrow b - Ax$ 
12:  else
13:     $r \leftarrow r - \alpha q$ 
14:  end if
15:   $\delta_{old} \leftarrow \delta_{new}$ 
16:   $\delta_{new} \leftarrow r^T r$ 
17:   $\beta \leftarrow \frac{\delta_{new}}{\delta_{old}}$ 
18:   $d \leftarrow r + \beta d$ 
19:   $i \leftarrow i + 1$ 
20: end while
```

PLSSVM - Parallel Least Squares Support Vector Machine

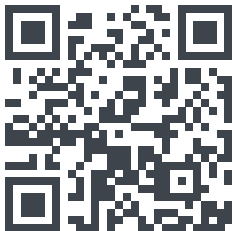
- modern C++17
- open source & on GitHub
- single and double precision via template parameter
- parallelizes matrix-vector multiplication in CG algorithm



<https://github.com/SC-SGS/PLSSVM>

PLSSVM - Parallel Least Squares Support Vector Machine

- modern C++17
- open source & on GitHub
- single and double precision via template parameter
- parallelizes matrix-vector multiplication in CG algorithm
- backends: **OpenMP, CUDA, HIP, OpenCL, and SYCL**
- backend and target platform selectable at runtime



<https://github.com/SC-SGS/PLSSVM>

PLSSVM - Parallel Least Squares Support Vector Machine

- modern C++17
- open source & on GitHub
- single and double precision via template parameter
- parallelizes matrix-vector multiplication in CG algorithm
- backends: **OpenMP, CUDA, HIP, OpenCL, and SYCL**
- backend and target platform selectable at runtime
- **multi-GPU** support for linear kernel function
- **drop-in replacement** for LIBSVM's `svm-train` and `svm-predict` executables
- currently only binary classification and dense calculations

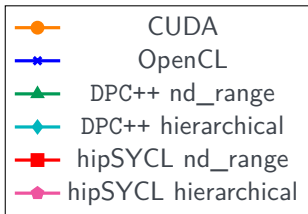
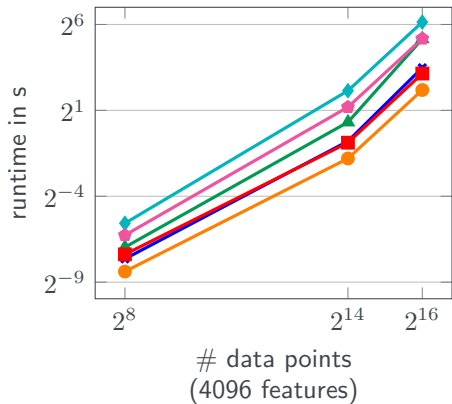


<https://github.com/SC-SGS/PLSSVM>

**New re-
sults and
findings**

2

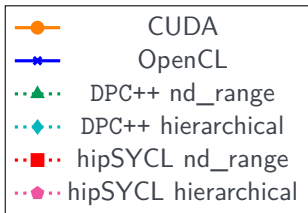
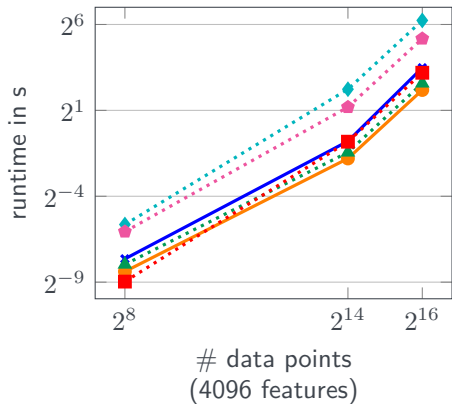
NVIDIA A100



Source: www.nvidia.com

	CUDA	OpenCL	DPC++ (20220202)		hipSYCL (Feb 1)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.003	0.005	0.008	0.021	0.006	0.013
16 384	0.287	0.576	1.242	4.418	0.543	2.281
65 536	4.547	11.113	35.709	70.418	8.848	36.102

NVIDIA A100



Source: www.nvidia.com

	CUDA	OpenCL	DPC++ (20221102)		hipSYCL (Oct 20)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.003	0.005	0.004 -50.0%	0.02 -4.76%	0.002 -66.66%	0.015 +15.38%
16 384	0.287	0.576	0.358 -71.18%	4.695 +6.27%	0.565 +3.89%	2.279 -0.1%
65 536	4.547	11.113	5.961 -83.31%	75.307 +6.94%	9.126 +3.05%	36.107 +0%

NVIDIA A100: explaining the results using profiling

16 384 x 4096	DPC++ 20220202	DPC++ 20221102	CUDA
runtime	1.242 s	0.358 s	0.287 s

NVIDIA A100: explaining the results using profiling

16 384 x 4096	DPC++ 20220202	DPC++ 20221102	CUDA
runtime	1.242 s	0.358 s	0.287 s
branch efficiency	65.06 %	99.97 %	99.97 %
avg divergent branches	3 972 456	170	170

NVIDIA A100: explaining the results using profiling

16 384 x 4096	DPC++ 20220202	DPC++ 20221102	CUDA
runtime	1.242 s	0.358 s	0.287 s
branch efficiency	65.06 %	99.97 %	99.97 %
avg divergent branches	3 972 456	170	170

sycl/handler.hpp (DPC++ 20220202)

```
1     template <typename KernelName, typename ElementType, typename KernelType>
2     __SYCL_KERNEL_ATTR__ void
3     #ifdef __SYCL_NONCONST_FUNCTOR__
4     kernel_parallel_for(KernelType KernelFunc) {
5     #else
6     kernel_parallel_for(const KernelType &KernelFunc) {
7     #endif
8     #ifdef __SYCL_DEVICE_ONLY__
9     KernelFunc(detail::Builder::getElement(detail::declptr<ElementType>()));
10    #else
11    (void)KernelFunc;
12    #endif
13    } // 925'493'095 (divergent branches)
```

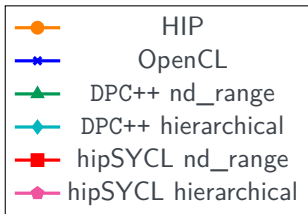
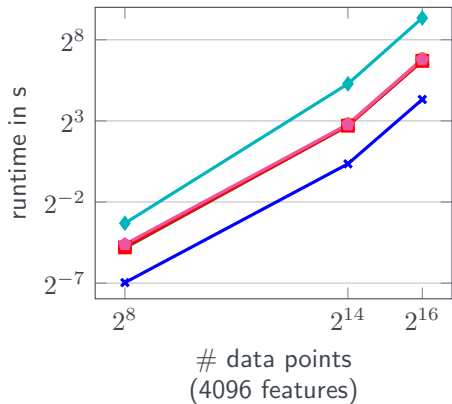
NVIDIA A100: explaining the results using profiling

16 384 x 4096	DPC++ 20220202	DPC++ 20221102	CUDA
runtime	1.242 s	0.358 s	0.287 s
branch efficiency	65.06 %	99.97 %	99.97 %
avg divergent branches	3 972 456	170	170
atomics (instr. exec.)	1 418 372 005	30 117 888	18 097 152

NVIDIA A100: explaining the results using profiling

16 384 x 4096	DPC++ 20220202	DPC++ 20221102	CUDA
runtime	1.242 s	0.358 s	0.287 s
branch efficiency	65.06 %	99.97 %	99.97 %
avg divergent branches	3 972 456	170	170
atomics (instr. exec.)	1 418 372 005	30 117 888	18 097 152
register count	164	164	162
memory	more memory transfers involving shared memory and between global \longleftrightarrow L1		better usage of registers; overall 43 % more memory throughput

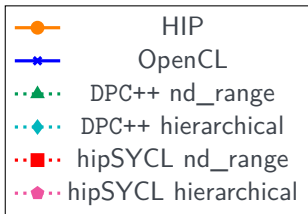
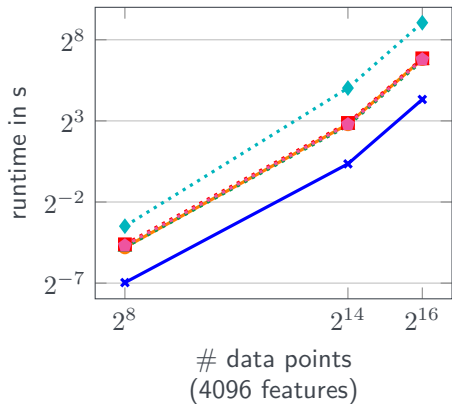
AMD Radeon Pro VII



Source: www.amd.com

	HIP	OpenCL	DPC++ (20220202)		hipSYCL (Feb 1)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.036	0.008	0.035	0.101	0.036	0.041
16 384	6.93	1.275	6.532	38.934	6.517	6.875
65 536	112.21	20.0	104.101	649.774	103.640	110.421

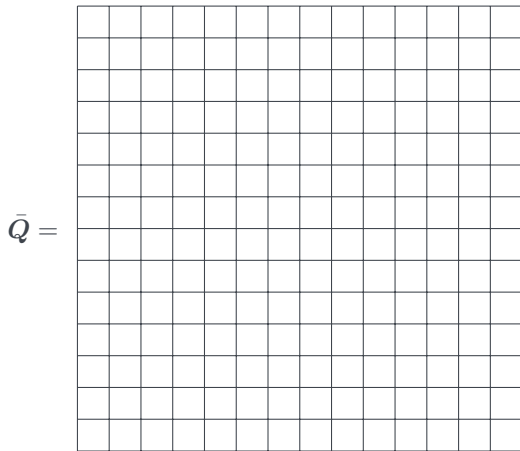
AMD Radeon Pro VII



Source: www.amd.com

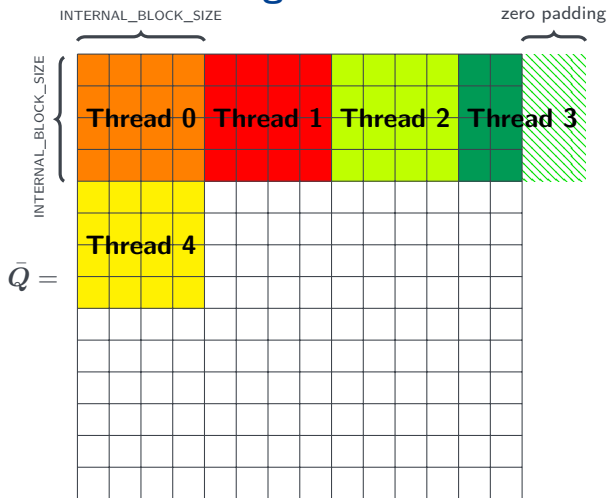
	HIP	OpenCL	DPC++ (20221102)		hipSYCL (Oct 20)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.036	0.008	0.036 +2.86 %	0.089 -11.88 %	0.041 +13.89 %	0.039 -4.88 %
16 384	6.93	1.275	6.547 +0.23 %	32.551 -16.39 %	7.327 +12.43 %	6.902 +0.39 %
65 536	112.21	20.0	104.366 +0.25 %	530.224 -18.40 %	115.807 +11.74 %	110.498 +0.08 %

Basic idea of the used blocking scheme



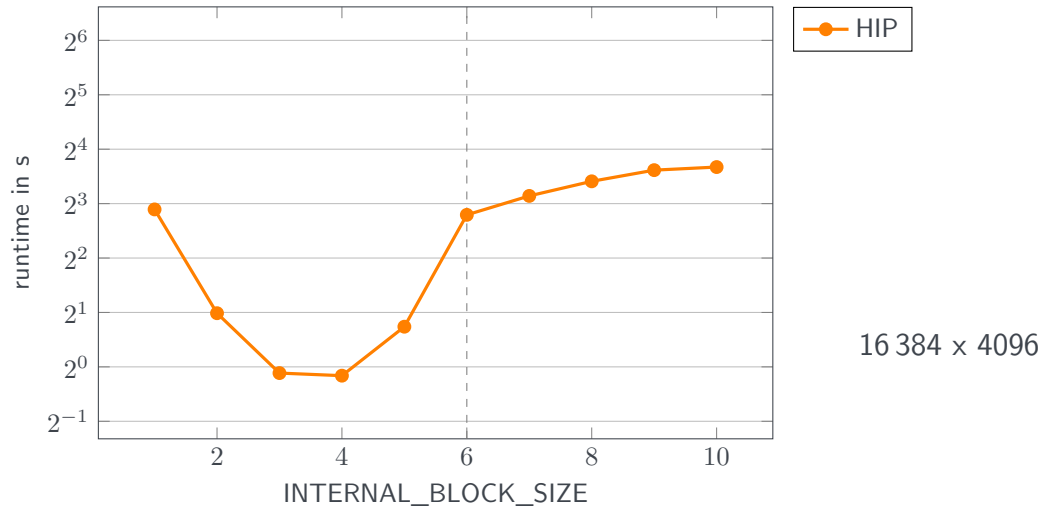
Note: each matrix entry Q_{ij} is calculated using the kernel function $k(\vec{x}_i, \vec{x}_j)$!
(e.g., dot products in the linear kernel)

Basic idea of the used blocking scheme

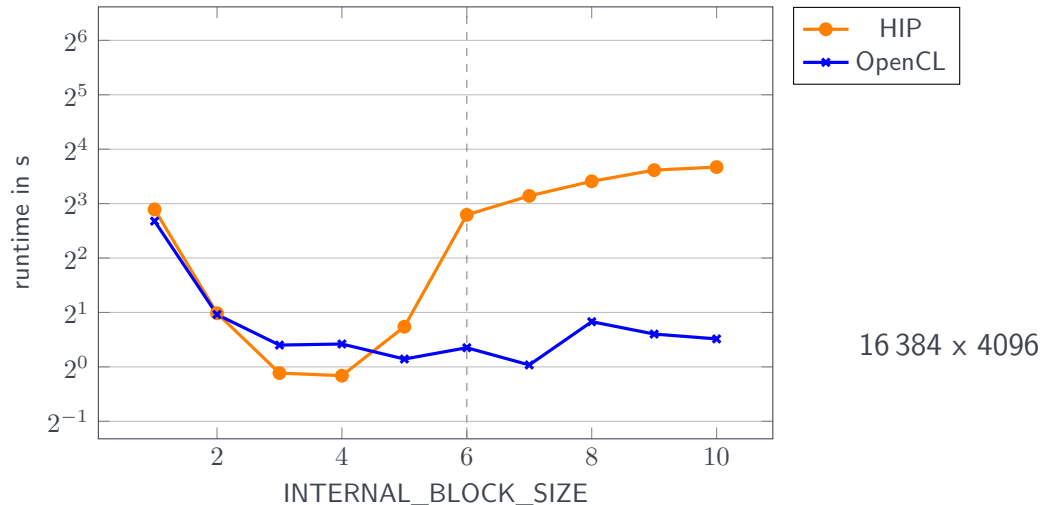


Note: each matrix entry Q_{ij} is calculated using the kernel function $k(\vec{x}_i, \vec{x}_j)$!
(e.g., dot products in the linear kernel)

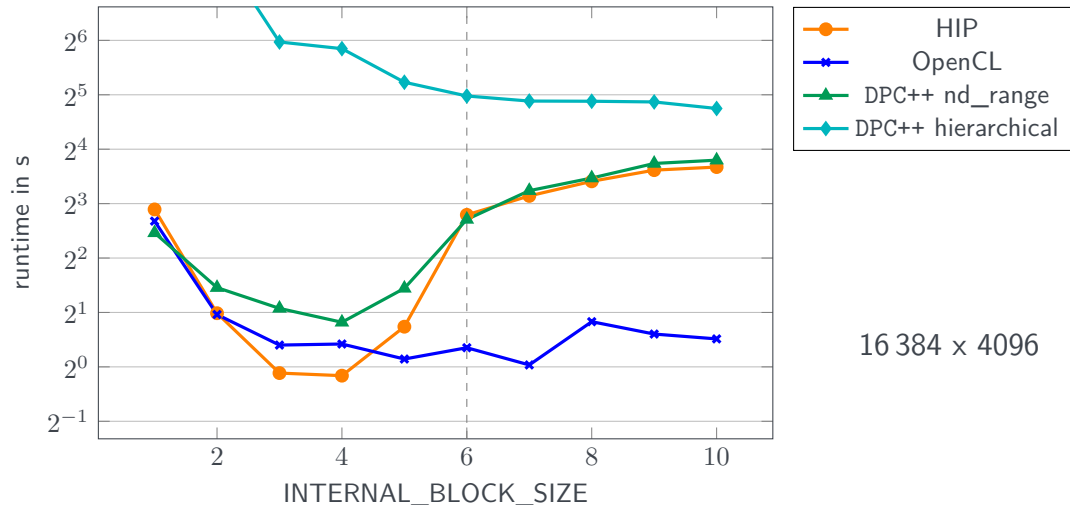
AMD Radeon Pro VII: Blocking Sizes



AMD Radeon Pro VII: Blocking Sizes

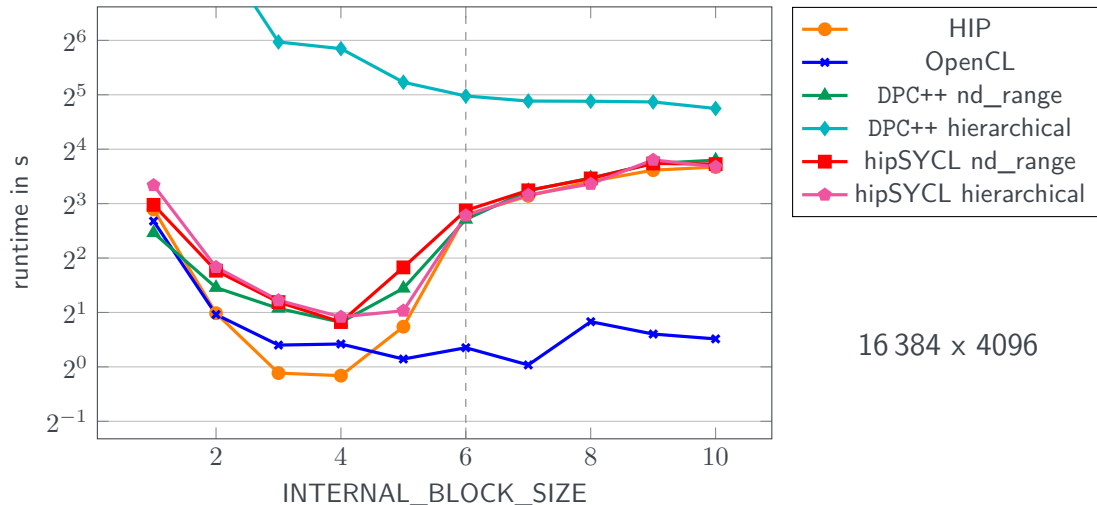


AMD Radeon Pro VII: Blocking Sizes

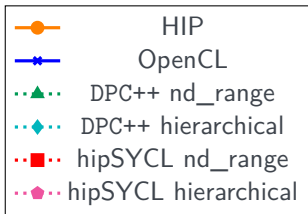
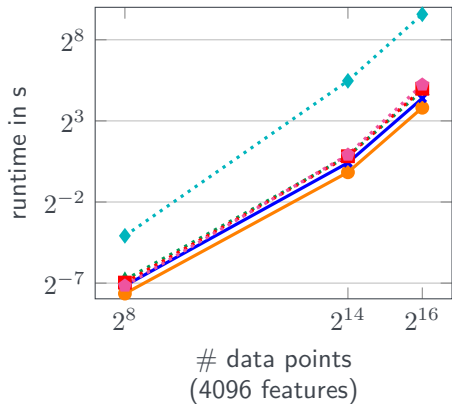


16384 x 4096

AMD Radeon Pro VII: Blocking Sizes



AMD Radeon Pro VII: updated runtimes with blocking size 4



Source: www.amd.com

	HIP	OpenCL	DPC++ (Nov 2)		hipSYCL (Oct 20)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.005 -86.11 %	0.007 -12.50 %	0.009 -75.00 %	0.059 -33.71 %	0.008 -80.49 %	0.007 -82.05 %
16 384	0.891 -87.14 %	1.335 +4.71 %	1.775 -72.89 %	44.243 -26.43 %	1.767 -75.88 %	1.882 -72.73 %
65 536	14.038 -87.49 %	20.995 +4.98 %	28.963 -72.25 %	761.951 +43.7 %	31.579 -72.73 %	37.568 -66.00 %

AMD Radeon Pro VII: explaining the results using profiling

16 384 x 4096	HIP		OpenCL	
INTERNAL_BLOCKING_SIZE	4	6	4	6
runtime	0.891 s	6.930 s	1.335 s	1.275 s

AMD Radeon Pro VII: explaining the results using profiling

16 384 x 4096	HIP		OpenCL	
INTERNAL_BLOCKING_SIZE	4	6	4	6
runtime	0.891 s	6.930 s	1.335 s	1.275 s
local data share	1024	1563	1024	1563
scratch memory	0	172	0	0
vector general purpose register	64	64	56	108

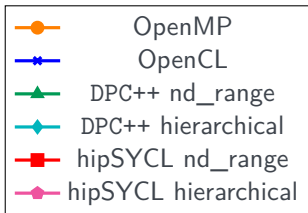
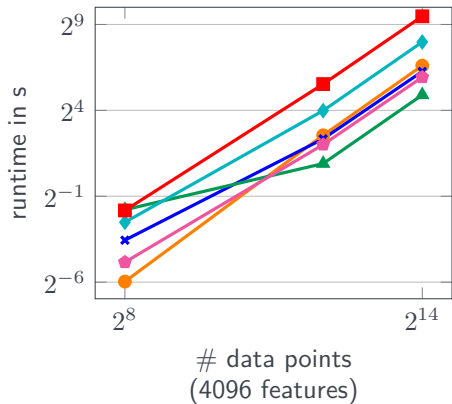
AMD Radeon Pro VII: explaining the results using profiling

16 384 x 4096	HIP		OpenCL	
INTERNAL_BLOCKING_SIZE	4	6	4	6
runtime	0.891 s	6.930 s	1.335 s	1.275 s
local data share	1024	1563	1024	1563
scratch memory	0	172	0	0
vector general purpose register	64	64	56	108
avg number of VALU inst.	57561	156393	58062	112165
% of active VALU threads	84.69 %	89.88 %	99.29 %	93.95 %

AMD Radeon Pro VII: explaining the results using profiling

16 384 x 4096	HIP		OpenCL	
INTERNAL_BLOCKING_SIZE	4	6	4	6
runtime	0.891 s	6.930 s	1.335 s	1.275 s
local data share	1024	1563	1024	1563
scratch memory	0	172	0	0
vector general purpose register	64	64	56	108
avg number of VALU inst.	57561	156393	58062	112165
% of active VALU threads	84.69 %	89.88 %	99.29 %	93.95 %
video memory fetches	84.29 GB	2039.79 GB	80.69 GB	53.48 GB
video memory writes	22.26 MB	1952.76 GB	19.45 MB	12.73 MB
bank conflicts (lower is better)	13.11 %	0.10 %	20.34 %	4.74 %

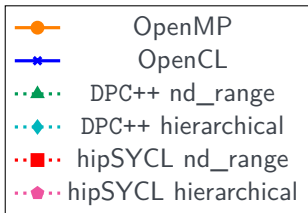
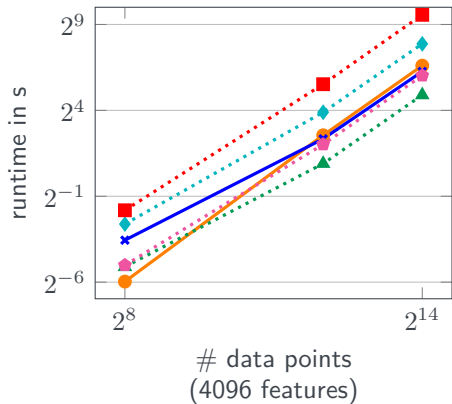
Intel Xeon E-2146G



Source: www.intel.com

	OpenMP	OpenCL	DPC++ (20220202)		hipSYCL (Feb 1)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.016	0.085	0.29	0.175	0.282	0.035
4096	5.855	5.066	1.869	15.819	46.201	4.020
16 384	97.161	76.765	29.844	252.245	711.613	61.042

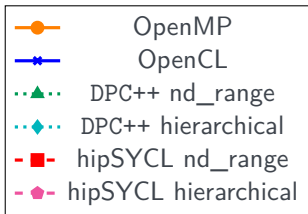
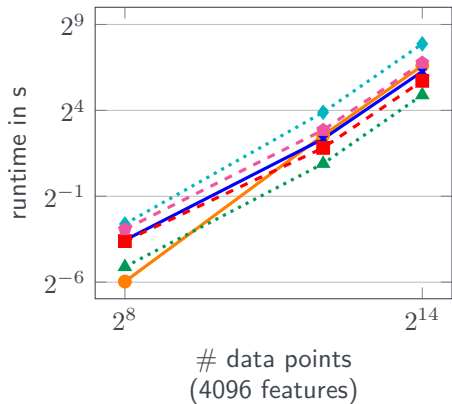
Intel Xeon E-2146G



Source: www.intel.com

	OpenMP	OpenCL	DPC++ (20221102)		hipSYCL (Oct 20)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.016	0.085	0.029 -90%	0.163 -6.86%	0.284 +0.71%	0.031 -11.43%
4096	5.855	5.066	1.866 -0.16%	14.806 -6.40%	45.954 -0.53%	4.049 +0.72%
16 384	97.161	76.765	29.730 -0.38%	234.304 -7.11%	755.112 +6.11%	65.149 +6.73%

Intel Xeon E-2146G



Source: www.intel.com

	OpenMP	OpenCL	DPC++ (20221102)		hipSYCL acc (Oct 20)	
			nd_range	hierarchical	nd_range	hierarchical
256	0.016	0.085	0.029 -90%	0.163 -6.86%	0.082 -70.92%	0.132 +277.14%
4096	5.855	5.066	1.866 -0.16%	14.806 -6.40%	3.521 -92.37%	7.235 +79.98%
16 384	97.161	76.765	29.730 -0.38%	234.304 -7.11%	52.715 -92.59%	109.161 +78.83%

Key takeaways: new versions improve the performance

	DPC++		hipSYCL	
	nd_range	hierarchical	nd_range	hierarchical
NVIDIA A100	↑	↘	→	→
AMD Radeon Pro VII	→	↑	↓	→
Intel Xeon E-2146G	→	↗	→/↑	→/↓

Key takeaways: the performance portability is good

Performance portability (application efficiency): (proposed by Pennycook, Sewall, and Lee in 2016)

$$\Phi(a, p, H) = \begin{cases} \frac{|H|}{\sum_{i \in H} \frac{1}{e_i(a, p)}} & \text{if } i \text{ is supported } \forall i \in H \\ 0 & \text{otherwise} \end{cases}$$

a : an application (modified matrix-vector multiplication)

p : a specific problem (16 384 × 4096)

H : a set of platforms (NVIDIA A100, AMD Radeon Pro VII, Intel Xeon)

Key takeaways: the performance portability is good

Performance portability (application efficiency): (proposed by Pennycook, Sewall, and Lee in 2016)

$$\mathbb{P}(a, p, H) = \begin{cases} \frac{|H|}{\sum_{i \in H} \frac{1}{e_i(a, p)}} & \text{if } i \text{ is supported } \forall i \in H \\ 0 & \text{otherwise} \end{cases}$$

a : an application (modified matrix-vector multiplication)

p : a specific problem (16 384 × 4096)

H : a set of platforms (NVIDIA A100, AMD Radeon Pro VII, Intel Xeon)

CUDA	HIP	OpenMP
0%	0%	0%

Key takeaways: the performance portability is good

Performance portability (application efficiency): (proposed by Pennycook, Sewall, and Lee in 2016)

$$\mathbb{P}(a, p, H) = \begin{cases} \frac{|H|}{\sum_{i \in H} \frac{1}{e_i(a, p)}} & \text{if } i \text{ is supported } \forall i \in H \\ 0 & \text{otherwise} \end{cases}$$

a : an application (modified matrix-vector multiplication)

p : a specific problem (16 384 × 4096)

H : a set of platforms (NVIDIA A100, AMD Radeon Pro VII, Intel Xeon)

CUDA	HIP	OpenMP	OpenCL	DPC++	hipSYCL
0 %	0 %	0 %	49.28 %	70.77 %	52.40 %

Conclusion

- fine-tuning hyper parameter (like the blocking size) can have a major impact on the performance
- profiling SYCL code (DPC++ and hipSYCL) is as easy as profiling native code

Conclusion

- fine-tuning hyper parameter (like the blocking size) can have a major impact on the performance
- profiling SYCL code (DPC++ and hipSYCL) is as easy as profiling native code
- installing new DPC++ or hipSYCL versions may drastically increase performance
- SYCL provides a better performance portability than OpenCL
 - in our case, DPC++ has the best performance portability with $\Phi(a, p, H) = 70.77\%$
- in addition: SYCL needs drastically less lines of code when compared to OpenCL
 - in our case, more the 300 lines of code

Conclusion

- fine-tuning hyper parameter (like the blocking size) can have a major impact on the performance
- profiling SYCL code (DPC++ and hipSYCL) is as easy as profiling native code
- installing new DPC++ or hipSYCL versions may drastically increase performance
- SYCL provides a better performance portability than OpenCL
 - in our case, DPC++ has the best performance portability with $\Phi(a, p, H) = 70.77\%$
- in addition: SYCL needs drastically less lines of code when compared to OpenCL
 - in our case, more the 300 lines of code

If performance portability is important, SYCL should be chosen over OpenCL!



University of Stuttgart
Germany



Marcel Breyer

Institute for Parallel and Distributed Systems
Scientific Computing



marcel.breyer@ipvs.uni-stuttgart.de



<https://orcid.org/0000-0003-3574-0650>

Further reading about PLSSVM

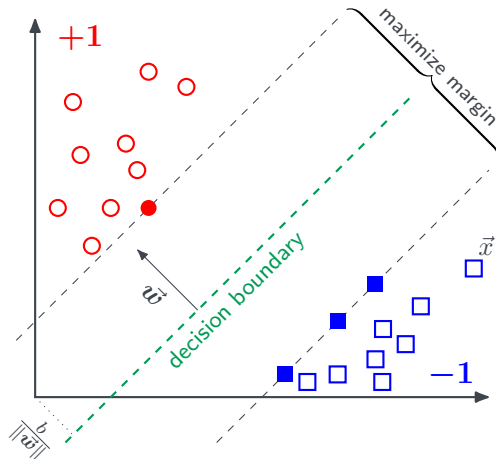
- [1] Alexander Van Craen, Marcel Breyer, and Dirk Pflüger. “PLSSVM: A (multi-)GPGPU-accelerated Least Squares Support Vector Machine”. In: *2022 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. 2022, pp. 818–827. DOI: [10.1109/IPDPSW55747.2022.00138](https://doi.org/10.1109/IPDPSW55747.2022.00138).
- [2] Marcel Breyer, Alexander Van Craen, and Dirk Pflüger. “A Comparison of SYCL, OpenCL, CUDA, and OpenMP for Massively Parallel Support Vector Machine Classification on Multi-Vendor Hardware”. In: *International Workshop on OpenCL. IWOCCL'22*. Bristol, United Kingdom, United Kingdom: Association for Computing Machinery, 2022. ISBN: 9781450396585. DOI: [10.1145/3529538.3529980](https://doi.org/10.1145/3529538.3529980). URL: <https://doi.org/10.1145/3529538.3529980>.
- [3] Alexander Van Craen, Marcel Breyer, and Dirk Pflüger. “PLSSVM—Parallel Least Squares Support Vector Machine”. In: *Software Impacts* 14 (2022), p. 100343. ISSN: 2665-9638. DOI: <https://doi.org/10.1016/j.simpa.2022.100343>. URL: <https://www.sciencedirect.com/science/article/pii/S2665963822000641>.



**Additional
resources**

Basics of Support Vector Machines (SVMs) (proposed by Boser, Guyon, and Vapnik in 1992)

supervised machine learning: binary classification



$$y = \text{sgn}(\langle \vec{w}, \vec{x} \rangle + b)$$

PLSSVM supports many different backends

OpenMP

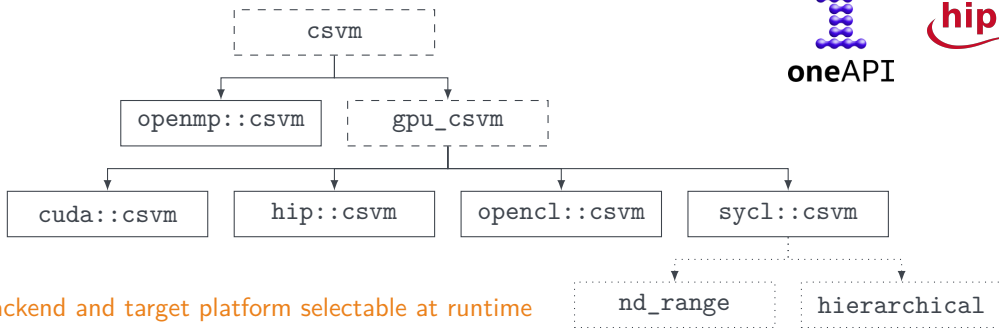


OpenCL

SYCL

oneAPI

hipSYCL



Different SYCL kernel invocation types

reverse all elements in an array

```
1  sycl::nd_range<1> exec{ global, local };
2  local_accessor<int> loc{ local , cgh };    // local memory
3  cgh.parallel_for(exec, [=](sycl::nd_item<1> item) {
4      const int idx = item.get_global_linear_id();
5      const int priv = n - idx - 1;        // private memory
6      loc[idx] = res[idx];
7      sycl::group_barrier(item.get_group()); // explicit barrier
8      res[idx] = loc[priv];
9  });
```

nd_range
(bottom-up)

(CUDA
HIP
OpenCL)

```
1  cgh.parallel_for_work_group(global, local, [=](sycl::group<1> group){
2      int loc[LOCAL_SIZE];                // local memory
3      sycl::private_memory<int> priv{ group }; // private memory
4      group.parallel_for_work_item([&](sycl::h_item<1> item) {
5          const int idx = item.get_local_id(0);
6          priv(item) = n - idx - 1;
7          loc[idx] = res[idx];
8      });
9      // implicit barrier
10     group.parallel_for_work_item([&](sycl::h_item<1> item) {
11         const int idx = item.get_local_id(0);
12         res[idx] = loc[priv(item)];
13     });
14 });
```

hierarchical
(top-down)

Used software and hardware



Source: www.nvidia.com



Source: www.amd.com



Source: www.intel.com

NVIDIA A100

CUDA 11.4.3

Driver Version 510.85.02

Radeon Pro VII

ROCm 5.3.0

Driver Version 5.18.2.22.40

Intel Xeon E-2146G

Intel DevCloud

DPC++

OpenSource LLVM fork

sycl-nightly/20220202 (February 02, 2022)

sycl-nightly/20221102 (November 02, 2022)

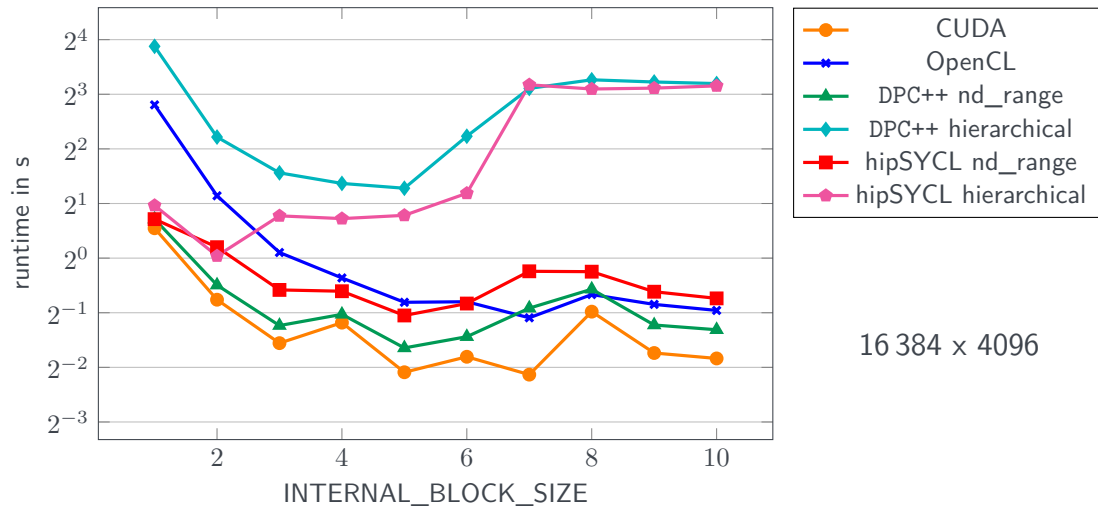
hipSYCL

OpenSource

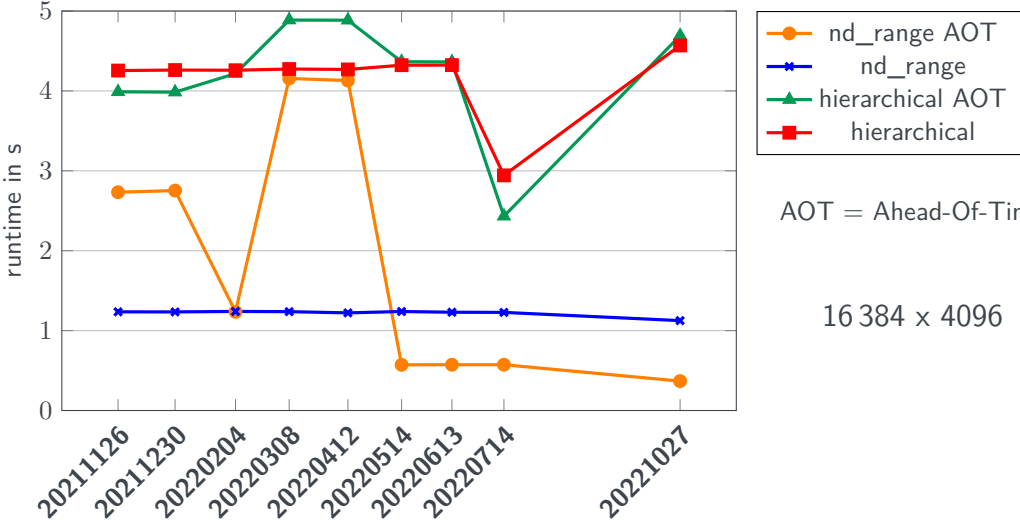
develop 6962942 (February 01, 2022)

develop 012e16d (October 20, 2022)

NVIDIA A100: varying blocking size



NVIDIA A100: the DPC++ compiler version makes a difference



AOT = Ahead-Of-Time

16 384 x 4096

Key takeaways: SYCL needs less lines of code than OpenCL

	kernel function	device discovery	other setup and bookkeeping code
CUDA	67	-	-
HIP	67	-	-
OpenMP	29	-	-
OpenCL	65	96	166 (kernel compilation & caching) 83 (custom sha256 for caching) 60 (3 custom RAII classes) 27 (custom atomic add) → 336
SYCL	nd_range		
	hierarchical	77	20 (used function object)