

IA³ 2020

10th Workshop on Irregular Applications: Architectures and Algorithms
November 11, 2020

Parallelizing Irregular Computations for Molecular Docking

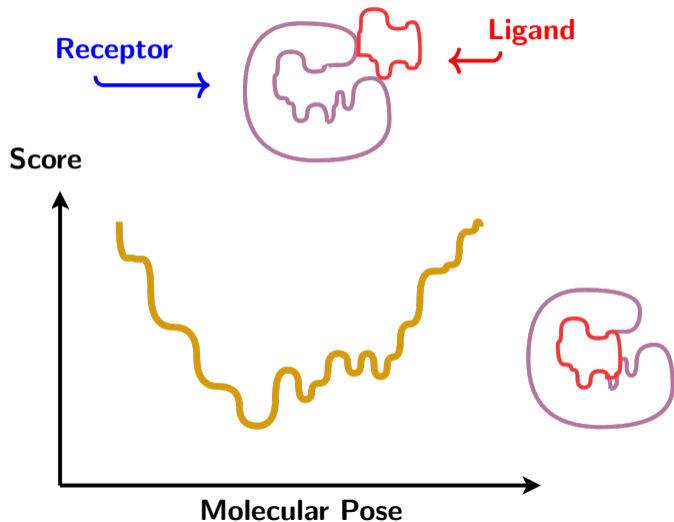
Leonardo Solis-Vasquez *, Diogo Santos-Martins ⁺, Andreas F. Tillack ⁺,

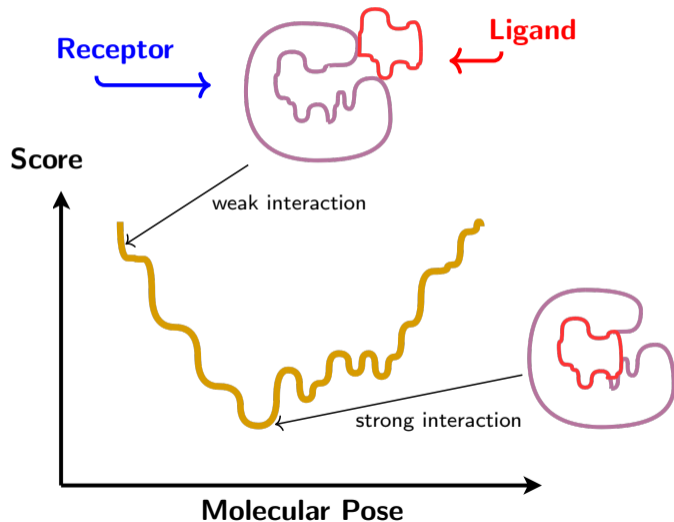
Andreas Koch *, Jérôme Eberhardt ⁺, Stefano Forli ⁺

* Embedded Systems and Applications Group
Technical University of Darmstadt, Germany

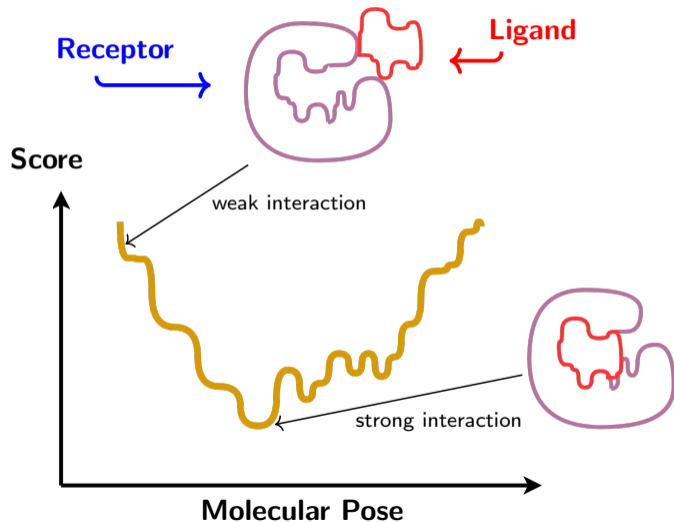
⁺ Center for Computational Structural Biology
The Scripps Research Institute, United States

Molecular Docking





- Molecular docking aims to find poses of strong interaction
- Scoring function
 - ▶ Measures how strong a pose is

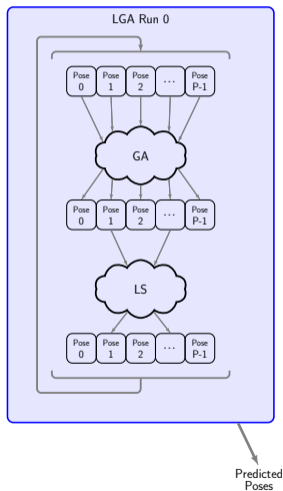


- Molecular docking aims to find poses of strong interaction
- Scoring function
 - ▶ Measures how strong a pose is
- Representation
 - ▶ Encodes a pose in terms of e.g., translation, rotation, torsion
- Search methods
 - ▶ Finds an optimal pose

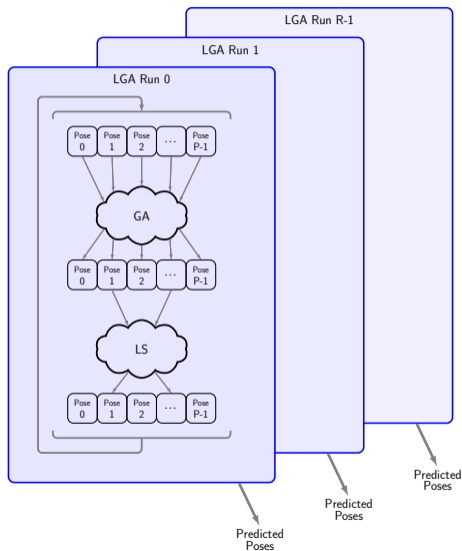
- One of the most cited docking tools
 - ▶ E.g., *FightAIDS@Home* project

- One of the most cited docking tools
 - ▶ E.g., *FightAIDS@Home* project
- Lamarckian Genetic Algorithm (LGA)
- LGA = GA + LS
 - ▶ Genetic Algorithm (GA)
 - ▶ Local Search (LS)
 - ▶ Both perform score calculations

- One of the most cited docking tools
 - ▶ E.g., *FightAIDS@Home* project
- Lamarckian Genetic Algorithm (LGA)
- LGA = GA + LS
 - ▶ Genetic Algorithm (GA)
 - ▶ Local Search (LS)
 - ▶ Both perform score calculations



- One of the most cited docking tools
 - ▶ E.g., *FightAIDS@Home* project
- Lamarckian Genetic Algorithm (LGA)
- LGA = GA + LS
 - ▶ Genetic Algorithm (GA)
 - ▶ Local Search (LS)
 - ▶ Both perform score calculations



- Binding energy (Kcal mol^{-1}) from molecular mechanics

- ▶ Molecular size

- ★ $N_{\text{atom}}^{\text{Receptor}} > 1000$

- ★ $N_{\text{atom}}^{\text{Ligand}} < 100$

- ▶ Dimensionless coefficients

- ★ $W_{\text{vdw}}, W_{\text{hb}}, W_{\text{el}}, W_{\text{ds}}, W_{\text{rot}}$

- ▶ Look-up tables

- ★ A, B, C, D, S, V, E, q

- ▶ Interatomic distance r_{ij}

- ★ Between atoms i and j

$$SF = \sum_{i,j} \left[\underbrace{W_{vdw} \left(\frac{A_{ij}}{r_{ij}^{12}} - \frac{B_{ij}}{r_{ij}^6} \right)}_{\text{Lennard-Jones}} + \underbrace{W_{hb} E(t) \left(\frac{C_{ij}}{r_{ij}^{12}} - \frac{D_{ij}}{r_{ij}^{10}} \right)}_{\text{Hydrogen bonding}} + \underbrace{W_{el} \left(\frac{q_i q_j}{\epsilon(r_{ij}) r_{ij}} \right)}_{\text{Coulomb's law}} + \underbrace{W_{ds} \left(S_i V_j + S_j V_i \right) e^{-\frac{r_{ij}^2}{2\sigma^2}}}_{\text{Desolvation}} \right]$$

- Termination criteria

- ▶ User defined

- ▶ $N_{\text{score-evals}}^{\text{MAX}} = 2\,048\,000$

- ▶ $N_{\text{gens}}^{\text{MAX}} = 27\,000$

- Nested loops

- ▶ With variable upper bounds
- ▶ Time-intensive score evals

```
Function AutoDock
```

```
  /* Coarse-Level Parallelism */  
  for each LGA-run do  
1     while ( $N_{\text{score-evals}} < N_{\text{score-evals}}^{\text{MAX}}$ ) and ( $N_{\text{gens}} < N_{\text{gens}}^{\text{MAX}}$ )  
2       do  
3         /* Medium-Level Parallelism */  
4         GA (population)  
5         /* Medium-Level Parallelism */  
         for individual in random-subset (population) do  
           LS (get-genotype (individual))
```

- Termination criteria

- ▶ User defined
- ▶ $N_{\text{LS-iters}}^{\text{MAX}} = 300$
- ▶ $\text{step}^{\text{MIN}} = 0,01$

- Nested loops

- ▶ With variable upper bounds
- ▶ Time-intensive score evals

- Divergent control

- ▶ Score improves \rightarrow *success*
- ▶ Score diminishes \rightarrow *failure*

```
/* Fine-Level Parallelism */
1 Function Solis-Wets (genotype)
2   while ( $N_{\text{LS-iters}} < N_{\text{LS-iters}}^{\text{MAX}}$ ) and (step >  $\text{step}^{\text{MIN}}$ ) do
3     delta = create-delta (step)
4     // new-genotype1
5     for each gene in  $N_{\text{genes}}$  do
6       new-gene1 = gene + delta
7     if SF (new-genotype1) < SF (genotype) then
8       genotype = new-genotype1
9       success++; fail = 0
10    else
11      // new-genotype2
12      for each gene in  $N_{\text{genes}}$  do
13        new-gene2 = gene - delta
14      if SF (new-genotype2) < SF (genotype) then
15        genotype = new-genotype2
16        success++; fail = 0
17      else
18        success = 0; fail++
19    step = update-step (success, fail)
```

- Previous work: OpenCL port of AutoDock
 - ▶ AutoDock-GPU
 - ★ Evaluated overall compute performance
 - ★ Focus: molecular prediction quality

- Previous work: OpenCL port of AutoDock
 - ▶ AutoDock-GPU
 - ★ Evaluated overall compute performance
 - ★ Focus: molecular prediction quality
- Here: parallelization in AutoDock-GPU
 - ▶ Focus: development rather than domain-oriented perspective
 - ▶ Challenges of dealing with AutoDock irregularity
 - ▶ Analysis of impact on execution runtime on GPUs/CPUs
 - ★ OpenCL work-groups configuration
 - ★ Molecular complexity of different inputs
 - ▶ Experiences porting onto FPGAs

Design Considerations for Host Code (1/2)

- AutoDock coded having only functionality in mind
 - ▶ I/O and compute tasks intertwine *unnecessarily*
 - ★ Read configuration options
 - ★ Perform computation (search and score calculation)
 - ★ Write partial results (predicted poses)
 - ★ Repeat (until all LGAs are processed)

Design Considerations for Host Code (1/2)

- AutoDock coded having only functionality in mind
 - ▶ I/O and compute tasks intertwine *unnecessarily*
 - ★ Read configuration options
 - ★ Perform computation (search and score calculation)
 - ★ Write partial results (predicted poses)
 - ★ Repeat (until all LGAs are processed)

- AutoDock-GPU re-structures program
 - ▶ Into a parallel-friendly version
 - ▶ I/O and compute tasks are decoupled *completely*
 - ▶ Exposes the Local Search function
 - ★ As the most runtime consuming
 - ★ Comprising several score evaluations

Design Considerations for Host Code (2/2)

- Rotatable bonds (torsions) affect ...
 - ▶ Interatomic distances → interactions
- AutoDock
 - ▶ Builds a tree of torsion-affected atoms
 - ▶ Recursively traverses that tree
 - ★ Calculates score at every node



Design Considerations for Host Code (2/2)

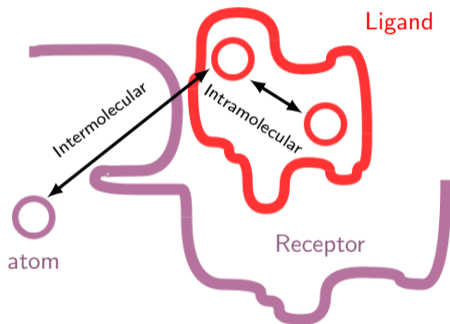
- Rotatable bonds (torsions) affect ...
 - ▶ Interatomic distances → interactions
- AutoDock
 - ▶ Builds a tree of torsion-affected atoms
 - ▶ Recursively traverses that tree
 - ★ Calculates score at every node
- AutoDock-GPU transforms data *on host*
 - ▶ Tree → arrays
 - ★ Recursion → iteration
 - ▶ More efficient on-device processing



ID	[Atom-to-Rotate Rotation-Type]
1	Rotation element 1
2	Rotation element 2
3	Rotation element 3
...	...
...	...
...	...
...	...
...	...
...	...
N	Rotation element N

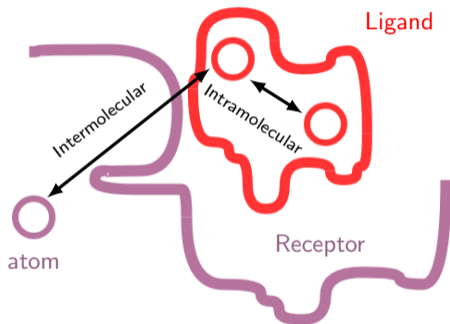
Re-designing Scoring Function (1/2)

- Scoring function has two components
 - ▶ Intermolecular
 - ★ receptor atoms \leftrightarrow ligand atoms
 - ▶ Intramolecular
 - ★ ligand atoms \leftrightarrow ligand atoms



Re-designing Scoring Function (1/2)

- Scoring function has two components
 - ▶ Intermolecular
 - ★ receptor atoms \leftrightarrow ligand atoms
 - ▶ Intramolecular
 - ★ ligand atoms \leftrightarrow ligand atoms
- AutoDock processes *pre-calculated* interactions
 - ▶ Purpose: reducing execution times
 - ▶ Pre-calculation takes place before AutoDock execution
 - ▶ Loop-up tables are accessed during docking



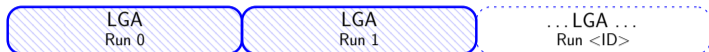
- AutoDock-GPU re-designs scoring
 - ▶ Pose calculation → integrated into scoring
 - ▶ Intermolecular → pre-calculated (still!)
 - ▶ Intramolecular
 - ★ Performs actual computations instead of pre-calculating
 - ★ $N_{\text{atom}}^{\text{Ligand}} < 100$ ($\ll N_{\text{atom}}^{\text{Receptor}}$)
 - ★ More accurate than pre-calculation
 - ★ Leverages compute power on e.g., GPUs

```
/* Fine-Level Parallelism */
1 Function SF (genotype)
2   for each rot-item in  $N_{\text{pose-rot}}$  do
3     | PoseCalculation
4   for each lig-atom in  $N_{\text{atom}}$  do
5     | InterInteraction
6   for each intra-pair in  $N_{\text{intra-contrib}}$  do
7     | IntraInteraction
```

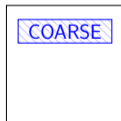
- AutoDock-GPU re-designs scoring
 - ▶ Pose calculation → integrated into scoring
 - ▶ Intermolecular → pre-calculated (still!)
 - ▶ Intramolecular
 - ★ Performs actual computations instead of pre-calculating
 - ★ $N_{\text{atom}}^{\text{Ligand}} < 100$ ($\ll N_{\text{atom}}^{\text{Receptor}}$)
 - ★ More accurate than pre-calculation
 - ★ Leverages compute power on e.g., GPUs
 - ▶ SF leverages fine-level parallelism

```
/* Fine-Level Parallelism */
1 Function SF (genotype)
2   for each rot-item in  $N_{\text{pose-rot}}$  do
3     | PoseCalculation
4   for each lig-atom in  $N_{\text{atom}}$  do
5     | InterInteraction
6   for each intra-pair in  $N_{\text{intra-contrib}}$  do
7     | IntraInteraction
```

Mapping operations into OpenCL elements



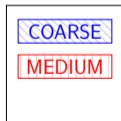
Parallelization
level



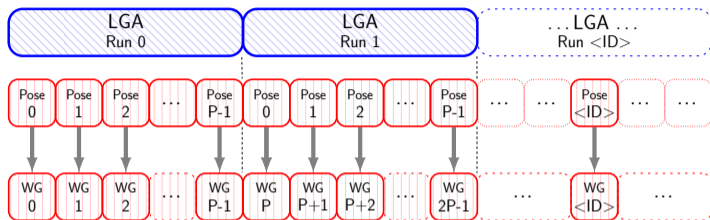
Mapping operations into OpenCL elements



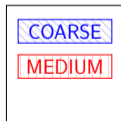
Parallelization
level



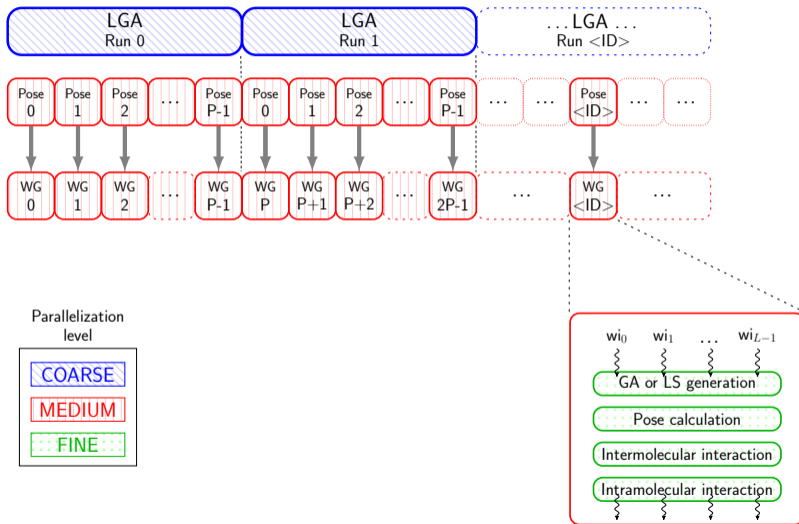
Mapping operations into OpenCL elements



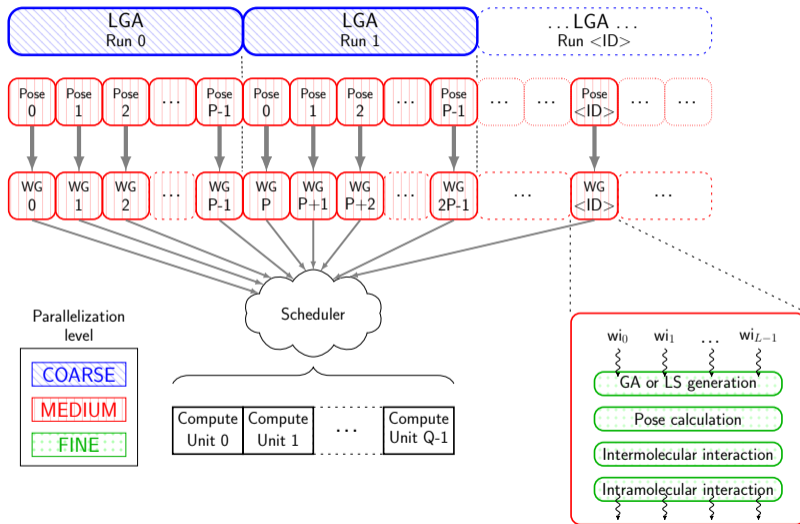
Parallelization
level



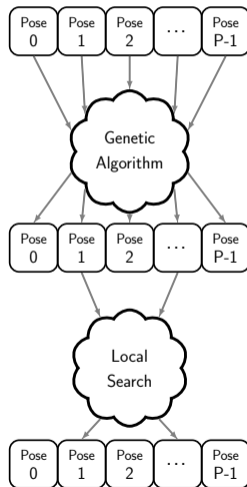
Mapping operations into OpenCL elements



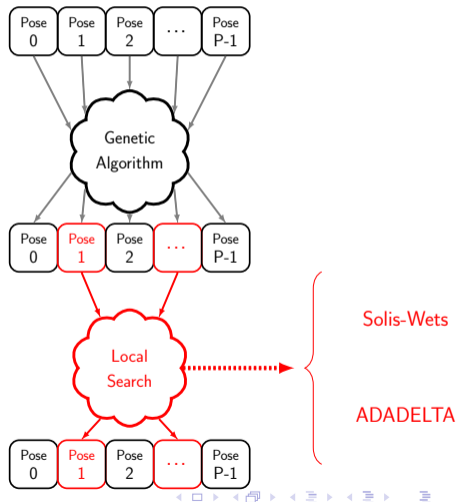
Mapping operations into OpenCL elements



- Local-search refinement
 - ▶ Can enhance pose predictions
- More efficient search algorithms
 - ▶ May find strong interactions
 - ▶ Require more-intensive computations



- Local-search refinement
 - ▶ Can enhance pose predictions
- More efficient search algorithms
 - ▶ May find strong interactions
 - ▶ Require more-intensive computations
- AutoDock-GPU code structure
 - ▶ Allows easy exchange between different local-search algorithms
 - ▶ Available local-search methods
 - ★ *Solis-Wets* (legacy)
 - ★ *ADADELTA* (newly incorporated)



```
/* Fine-Level Parallelism */
1 Function GC (genotype)
  /* Gradients in atomic space */
2   for each rot-item in  $N_{\text{pose-rot}}$  do
3     | PoseCalculation
4   for each lig-atom in  $N_{\text{atom}}$  do
5     | InterGradient
6   for each intra-pair in  $N_{\text{intra-contrib}}$  do
7     | IntraGradient
  /* Conversion into genetic space */
8   Gtrans // Translational gradients
9   Grigidrot // Rigid-body rotation gradients
10  Grotbond // Rotatable-bond gradients
```

```
/* Fine-Level Parallelism */
1 Function GC (genotype)
  /* Gradients in atomic space */
2   for each rot-item in  $N_{\text{pose-rot}}$  do
3     | PoseCalculation
4   for each lig-atom in  $N_{\text{atom}}$  do
5     | InterGradient
6   for each intra-pair in  $N_{\text{intra-contrib}}$  do
7     | IntraGradient
  /* Conversion into genetic space */
8   Gtrans // Translational gradients
9   Grigidrot // Rigid-body rotation gradients
10  Grotbond // Rotatable-bond gradients
```

```
/* Fine-Level Parallelism */
1 Function ADADELTA (genotype)
2   gradient = GC (genotype)
3   while ( $N_{\text{LS-iters}} < N_{\text{LS-iters}}^{\text{MAX}}$ ) do
4     | new-genotype = update-rule (genotype, gradient)
5     | if SF (new-genotype) < SF (genotype) then
6       | | genotype = new-genotype
7     | | gradient = GC (genotype)
```

- Parameters

- ▶ NDR_{size} : # work-items per kernel
- ▶ R : # LGA runs
- ▶ P : population size
- ▶ $lsrate$: local-search rate
- ▶ WG_{size} : # work-items per work-group

$$NDR_{\text{size}}^{\text{KrnL-GA}} = \{R \times P \times WG_{\text{size}}, 1, 1\}$$

$$NDR_{\text{size}}^{\text{KrnL-LS}} = \{R \times P \times lsrate \times WG_{\text{size}}, 1, 1\}$$

- Parameters

- ▶ NDR_{size} : # work-items per kernel
- ▶ R : # LGA runs
- ▶ P : population size
- ▶ $l\text{rate}$: local-search rate
- ▶ WG_{size} : # work-items per work-group

$$NDR_{\text{size}}^{\text{KrnL-GA}} = \{R \times P \times WG_{\text{size}}, 1, 1\}$$

$$NDR_{\text{size}}^{\text{KrnL-LS}} = \{R \times P \times l\text{rate} \times WG_{\text{size}}, 1, 1\}$$

- For all experiments

- ▶ $R = 100$
- ▶ $P = 150$
- ▶ $l\text{rate} = 100\%$
- If $WG_{\text{size}}^{\text{GPU}} = 64 \rightarrow NDR_{\text{size}}^{\text{GPU}} = \{960000, 1, 1\}$
- If $WG_{\text{size}}^{\text{CPU}} = 16 \rightarrow NDR_{\text{size}}^{\text{CPU}} = \{240000, 1, 1\}$

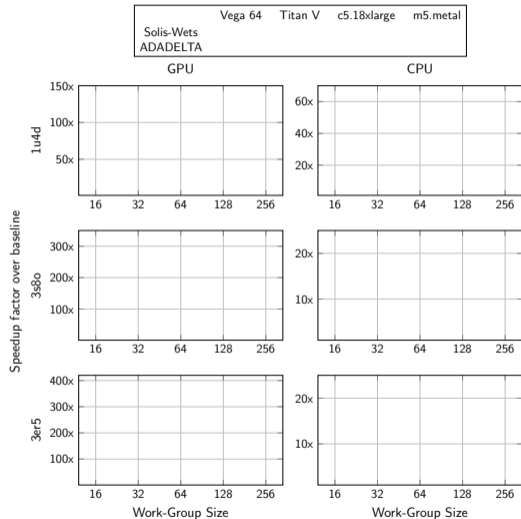
- AutoDock: v4.2.6 (baseline)
 - ▶ Implements only Solis-Wets method
 - ▶ Does not support multithreading
 - ★ Run on a Xeon Platinum 8124M @3.0 GHz CPU core

- AutoDock: v4.2.6 (baseline)
 - ▶ Implements only Solis-Wets method
 - ▶ Does not support multithreading
 - ★ Run on a Xeon Platinum 8124M @3.0 GHz CPU core

- AutoDock-GPU: v1.2
 - ▶ Implements both Solis-Wets and ADADELTA methods
 - ▶ Run on different GPU/CPU accelerators
 - ★ Radeon RX Vega 64 GPU (on-premise)
 - ★ Volta Titan V GPU (on-premise)
 - ★ Xeon Platinum 8124M @3.0 GHz 36-core CPU (AWS c5.18xlarge)
 - ★ Xeon Platinum 8175M @2.5 GHz 48-core CPU (AWS m5.meta1)

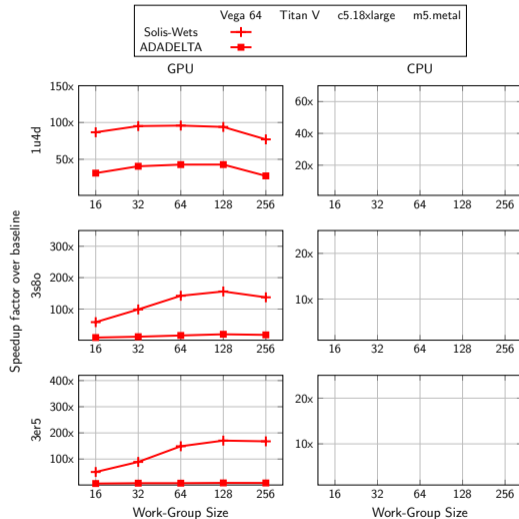
Speedup vs. OpenCL Work-Group Size

- Inputs of different complexity
 - ▶ 1u4d (low), 3s8o (medium), 3er5 (high)



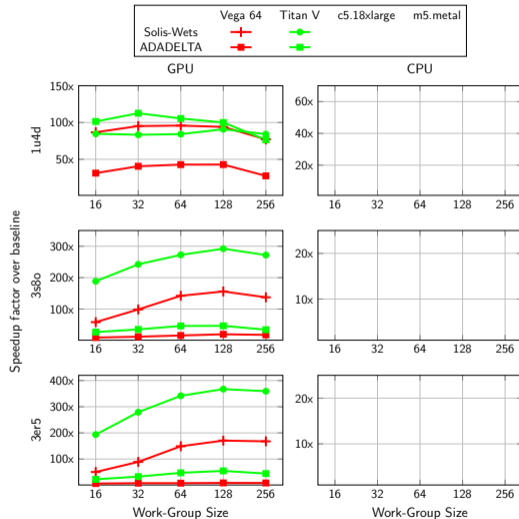
Speedup vs. OpenCL Work-Group Size

- Inputs of different complexity
 - ▶ 1u4d (low), 3s8o (medium), 3er5 (high)
- GPUs: best WG_{size} depends on ...
 - ▶ Molecular complexity
 - ▶ Accelerator being employed



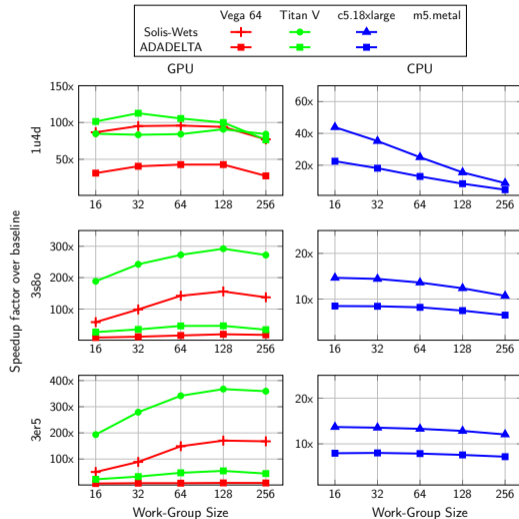
Speedup vs. OpenCL Work-Group Size

- Inputs of different complexity
 - ▶ 1u4d (low), 3s8o (medium), 3er5 (high)
- GPUs: best WG_{size} depends on ...
 - ▶ Molecular complexity
 - ▶ Accelerator being employed



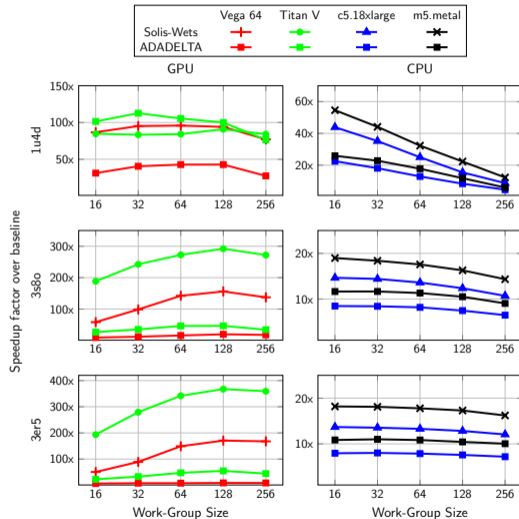
Speedup vs. OpenCL Work-Group Size

- Inputs of different complexity
 - ▶ 1u4d (low), 3s8o (medium), 3er5 (high)
- GPUs: best WG_{size} depends on ...
 - ▶ Molecular complexity
 - ▶ Accelerator being employed
- CPUs: faster executions when ...
 - ▶ Smaller WG_{size}
 - ▶ $WG_{size} = 16$



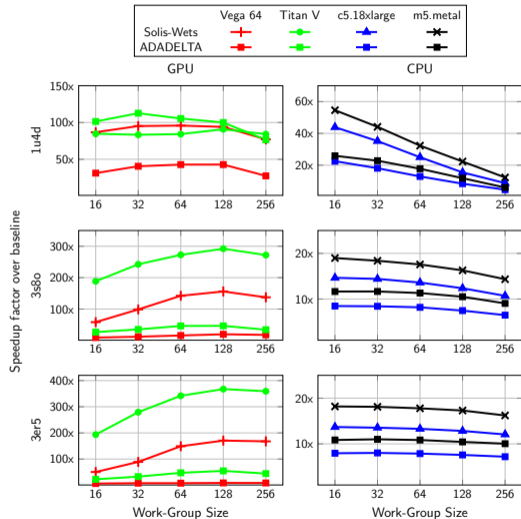
Speedup vs. OpenCL Work-Group Size

- Inputs of different complexity
 - ▶ 1u4d (low), 3s8o (medium), 3er5 (high)
- GPUs: best WG_{size} depends on ...
 - ▶ Molecular complexity
 - ▶ Accelerator being employed
- CPUs: faster executions when ...
 - ▶ Smaller WG_{size}
 - ▶ $WG_{size} = 16$

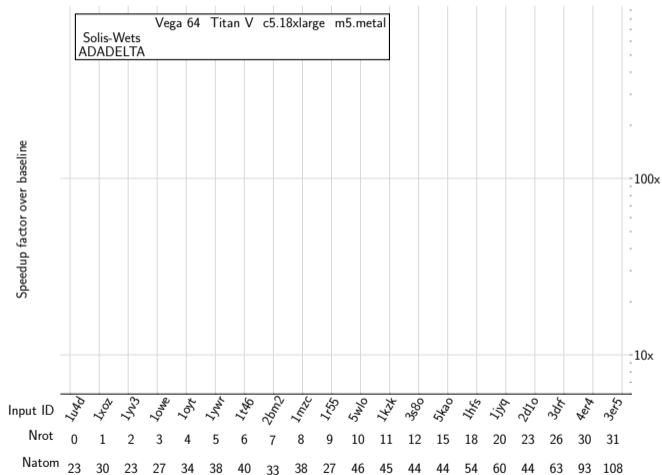


Speedup vs. OpenCL Work-Group Size

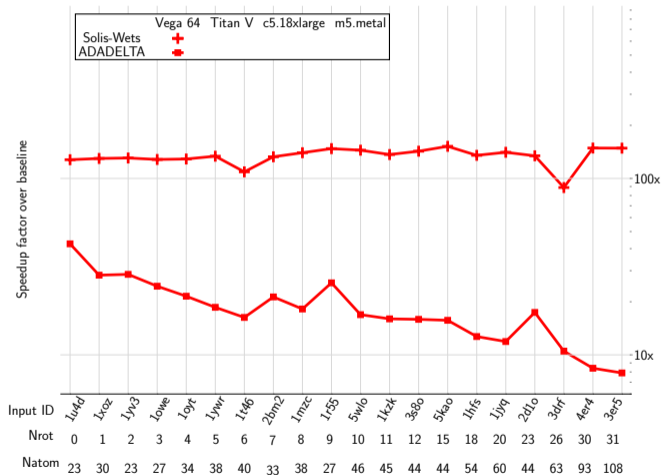
- Inputs of different complexity
 - ▶ 1u4d (low), 3s8o (medium), 3er5 (high)
- GPUs: best WG_{size} depends on ...
 - ▶ Molecular complexity
 - ▶ Accelerator being employed
- CPUs: faster executions when ...
 - ▶ Smaller WG_{size}
 - ▶ $WG_{size} = 16$
- For next experiments
 - ▶ $WG_{size}^{CPU} = 16$
 - ▶ $WG_{size}^{GPU} = 64$
 - ★ Min. multiple of a Nvidia *warp* (32) and AMD *wavefront* (64)



Speedup vs. Molecular Complexity

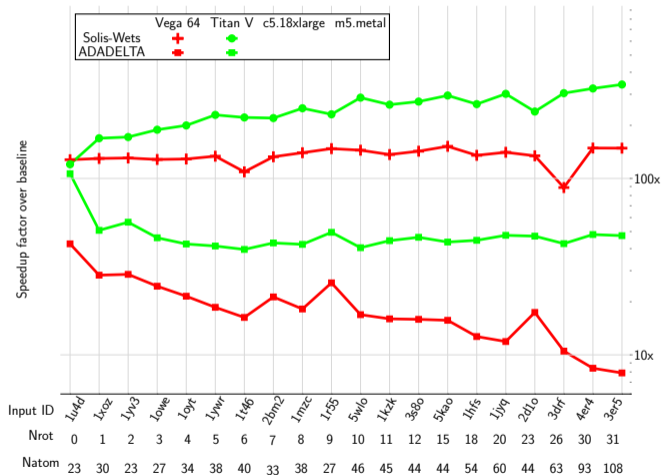


Speedup vs. Molecular Complexity

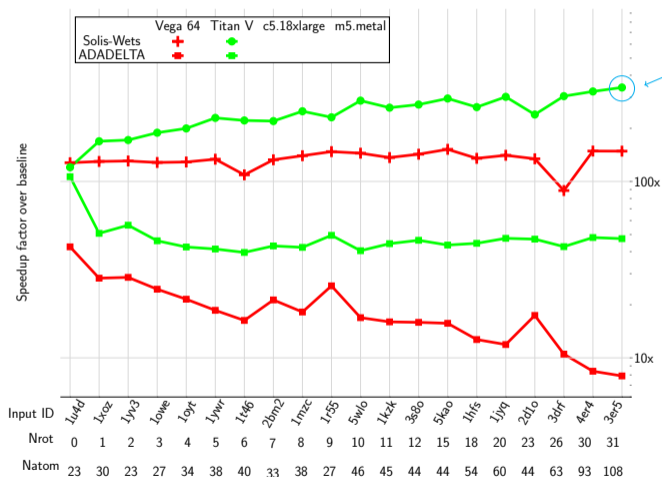


Growing molecular complexity →

Speedup vs. Molecular Complexity



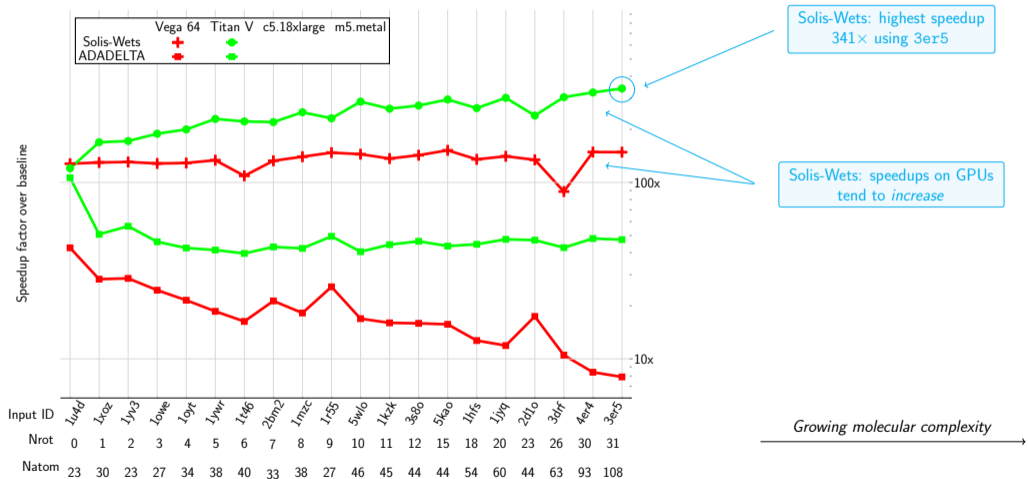
Speedup vs. Molecular Complexity



Solis-Wets: highest speedup
341x using 3e5

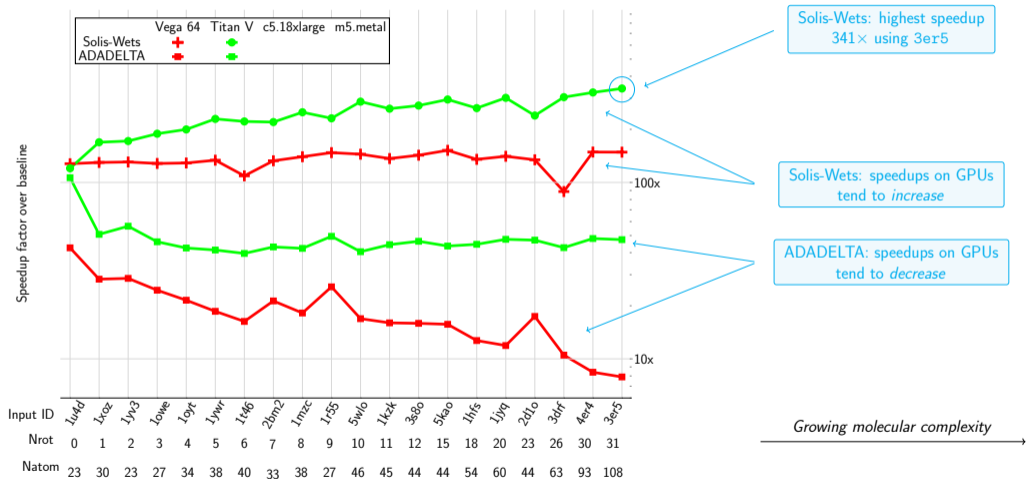
Growing molecular complexity →

Speedup vs. Molecular Complexity

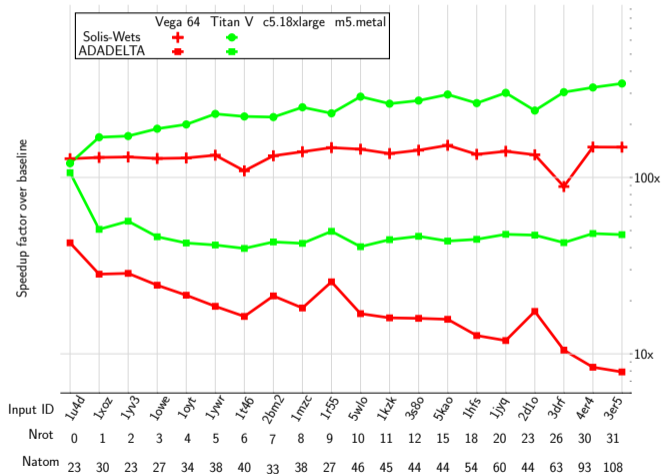


Growing molecular complexity →

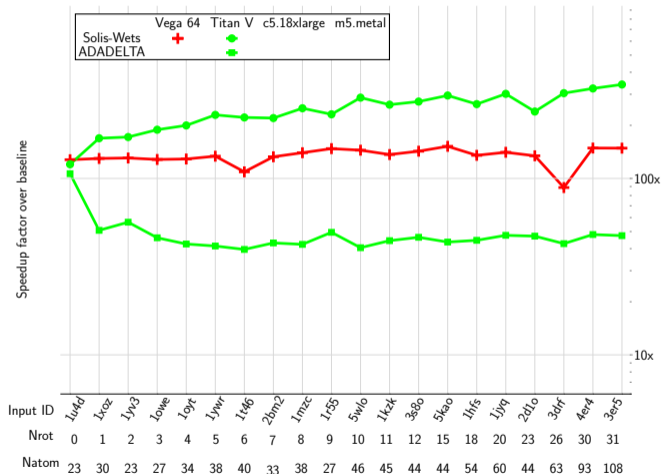
Speedup vs. Molecular Complexity



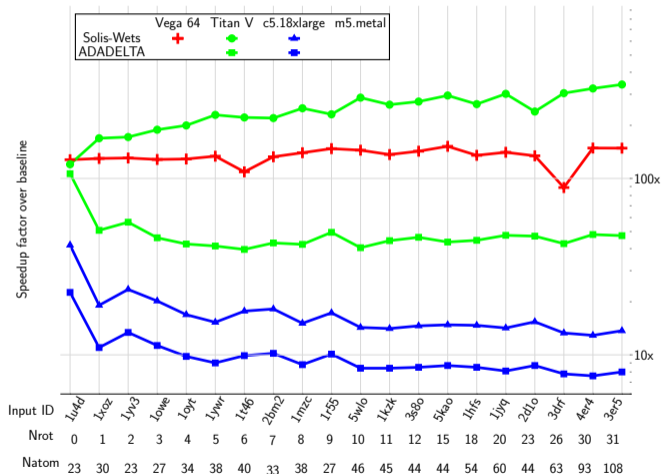
Speedup vs. Molecular Complexity



Speedup vs. Molecular Complexity

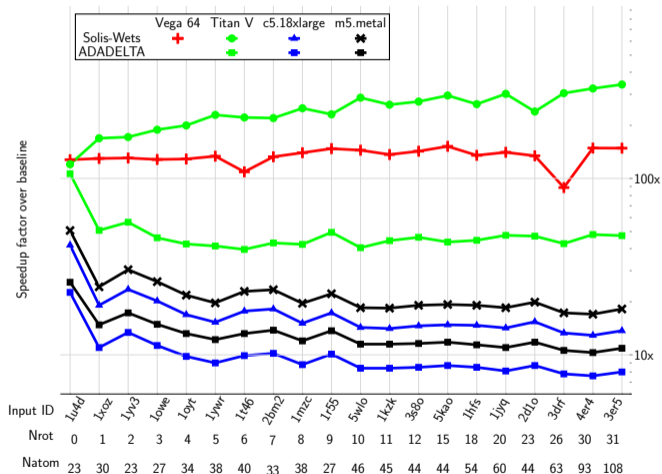


Speedup vs. Molecular Complexity



Growing molecular complexity →

Speedup vs. Molecular Complexity



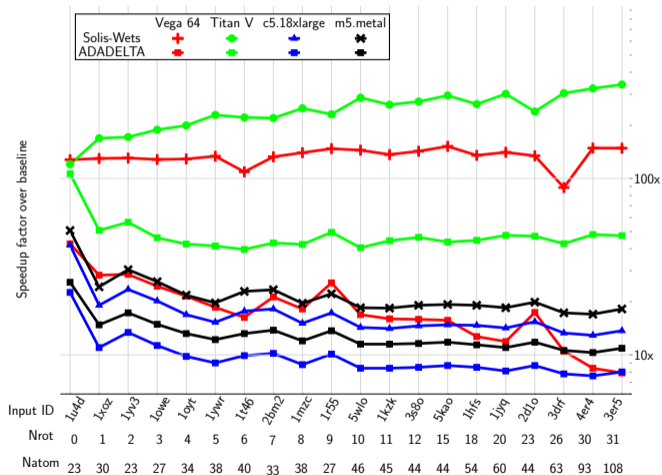
Solis-Wets & ADADELTA:
speedups on CPUs tend to *decrease*

Similar profiles among CPUs

- m5.metal: 48 physical cores
- c5.18xlarge: 36 physical cores
- Speedup ratio: $<1.2\times - 1.4\times>$

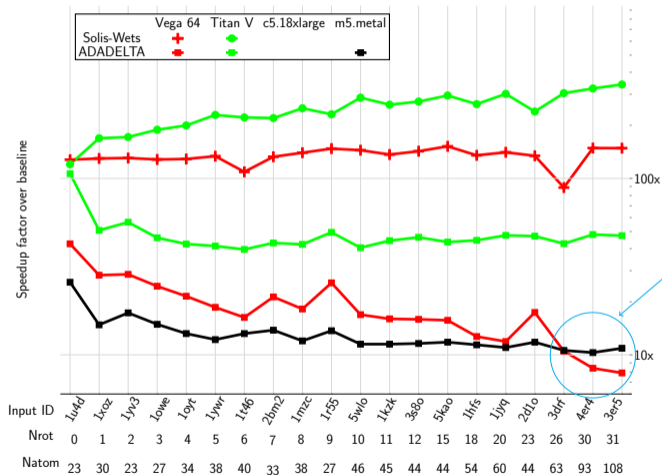
Growing molecular complexity →

Speedup vs. Molecular Complexity



Growing molecular complexity →

Speedup vs. Molecular Complexity

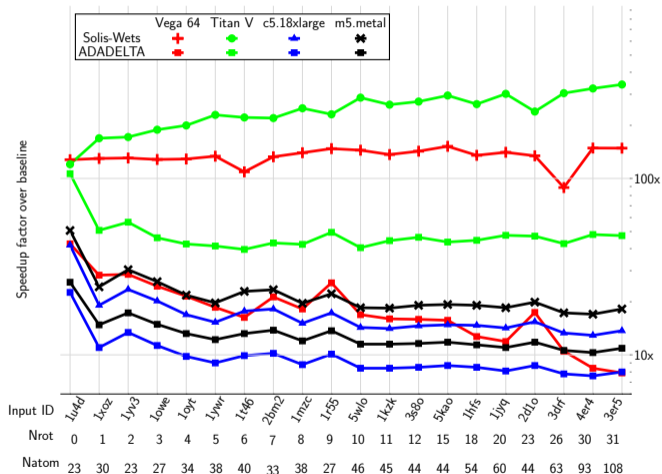


GPUs achieve higher speedups than CPUs
(in most cases)

Exception: 3drf, 4er4, 3er5
 - Vega 64 GPU: 11x, 8x, 8x
 - m5.metal CPU: 11x, 10x, 11x

Growing molecular complexity →

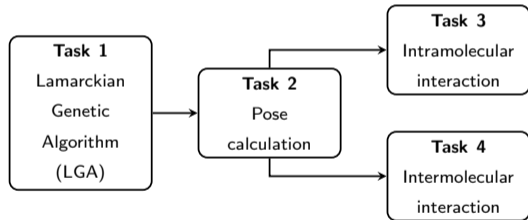
Speedup vs. Molecular Complexity



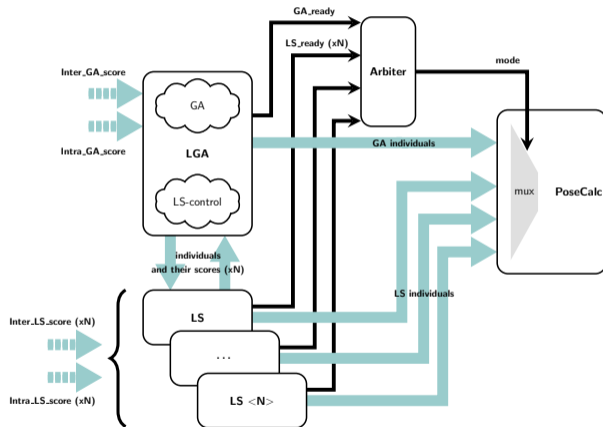
Efficiency:
Solis-Wets > ADADELTA

Growing molecular complexity →

- Data-parallel design on FPGAs
 - ▶ Three orders of magnitude *slower*
- Task parallelization
 - ▶ Each task coded as a single work-item kernel
 - ▶ Kernels communicate via OpenCL pipes
 - ▶ General design practices
 - ★ Pipelining loops within each kernel
 - ★ Minimizing loops initiation interval



- Final design composed of 27 kernels
 - ▶ Additional kernels (and pipes)
 - ★ Local-search kernels (Solis-Wets)
 - ★ Random number generators



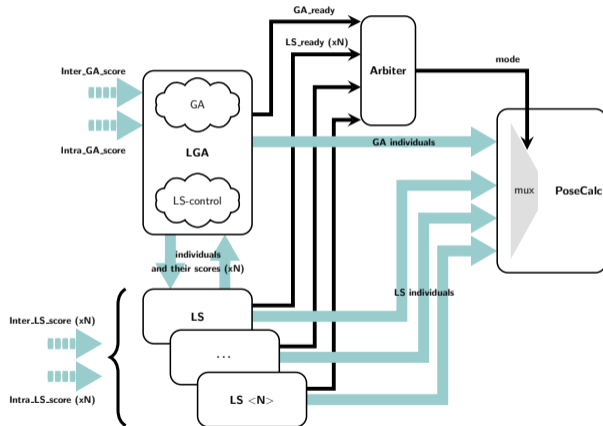
- Final design composed of 27 kernels

- ▶ Additional kernels (and pipes)

- ★ Local-search kernels (Solis-Wets)
- ★ Random number generators

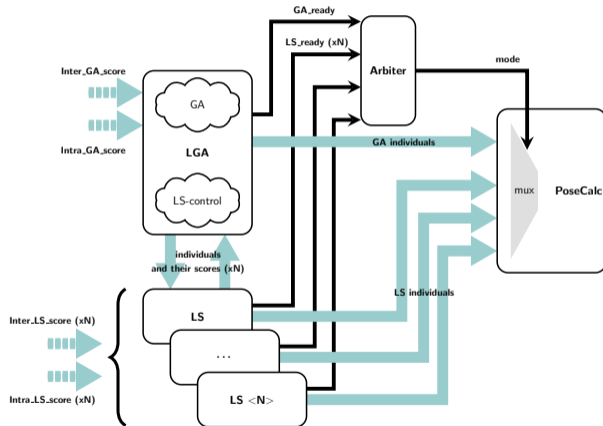
- ▶ Complex synchronization

- ★ All kernels running *simultaneously*
- ★ Pipes configured as *blocking* and *non-blocking*
- ★ Non-blocking pipes handle variable kernel communication



- Final design composed of 27 kernels

- ▶ Additional kernels (and pipes)
 - ★ Local-search kernels (Solis-Wets)
 - ★ Random number generators
- ▶ Complex synchronization
 - ★ All kernels running *simultaneously*
 - ★ Pipes configured as *blocking* and *non-blocking*
 - ★ Non-blocking pipes handle variable kernel communication
- ▶ Lower speedups on FPGAs wrt. GPUs
 - ★ Arria 10: 3× faster than serial baseline

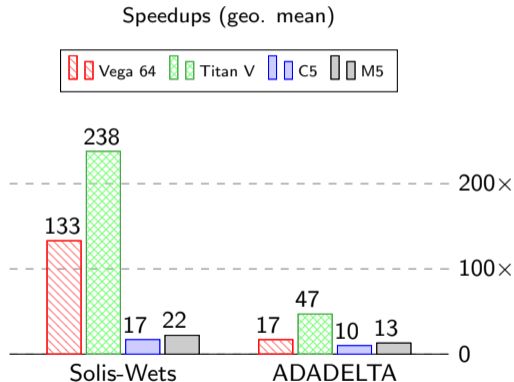


- AWS f1.2xlarge instance
 - ▶ Compilation successful after minor code changes
 - ▶ Execution on FPGA fails
 - ★ Non-blocking pipes are *not supported* in Xilinx tools
 - ▶ Possible avenue
 - ★ E.g., replacing variable by constant upper-bounds for loops
 - ★ Using only blocking pipes (supported!)
 - ▶ Transforming into a *regular* application ?

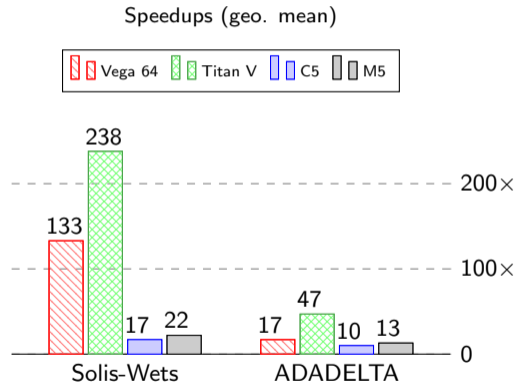
- Parallelized AutoDock using OpenCL
 - ▶ Challenges due to irregularity in AutoDock
 - ★ Divergent control performing local search
 - ★ Loops with variable upper bounds
 - ★ Time-intensive score evaluations
 - ▶ Required large-scale code re-structuring
 - ★ Trees → arrays
 - ★ Score function re-structuring

- Parallelized AutoDock using OpenCL
 - ▶ Challenges due to irregularity in AutoDock
 - ★ Divergent control performing local search
 - ★ Loops with variable upper bounds
 - ★ Time-intensive score evaluations
 - ▶ Required large-scale code re-structuring
 - ★ Trees → arrays
 - ★ Score function re-structuring
 - ▶ OpenCL work-items
 - ★ Atomic rotations and score calculations
 - ▶ OpenCL work-groups
 - ★ Simultaneously processing molecular poses

- Evaluated AutoDock-GPU performance
 - ▶ Comparing local search
 - ★ Solis-Wets vs. ADADELTA
 - ★ Solis-Wets: higher speedups
 - ★ ADADELTA: better pose predictions
 - ▶ Overall performance
 - ★ Depends on the input molecule



- Evaluated AutoDock-GPU performance
 - ▶ Comparing local search
 - ★ Solis-Wets vs. ADADELTA
 - ★ Solis-Wets: higher speedups
 - ★ ADADELTA: better pose predictions
 - ▶ Overall performance
 - ★ Depends on the input molecule
- Lower speedups achieved on FPGAs
 - ▶ Due to irregularity in AutoDock



Parallelizing Irregular Computations for Molecular Docking

<https://github.com/ccsb-scripps/AutoDock-GPU>

<https://www.esa.tu-darmstadt.de>

<https://forlilab.org>