

Parallelizing with xDSC, a Resource-Constrained Scheduling Algorithm for Shared and Distributed Memory Systems

Dounia KHALDI

CRI, Mathématiques et systèmes
MINES ParisTech

June 19, 2012

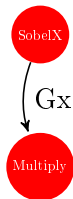


Evolution of the architecture (Multicores, GPUs...)

Evolution of parallel execution environments (OpenMP, MPI, OpenCL...)

Parallel software developed by converting sequential programs by hand

```
in = InitHarris();
//Sobel
SobelX(Gx, in);
SobelY(Gy, in);
//Multiply
MultiplY(Ixx, Gx, Gx);
MultiplY(Iyy, Gy, Gy);
MultiplY(Ixy, Gx, Gy);
//Gauss
Gauss(Sxx, Ixx);
Gauss(Syy, Iyy);
Gauss(Sxy, Ixy);
//Coarsity
CoarsitY(out, Sxx, Syy, Sxy);
```



- Scheduling \Rightarrow minimize completion time.
- $\text{length}(\text{path}) = \text{communication_cost}(\text{edges}) + \text{computational_cost}(\text{nodes})$.
- Dynamic vs. Static.
- List-scheduling heuristics.

List-Scheduling Processes

- Priorities are computed of all unscheduled nodes:
 - Top level ($tlevel(\tau)$): length of the longest path from the entry node to $\tau \Rightarrow$ earliest possible start-time.
 - Bottom level ($blevel(\tau)$): length of the longest path from τ to the exit node \Rightarrow latest start-time = $CriticalPathLength - blevel(\tau)$.
- The node τ with the highest priority is selected for scheduling.
- τ is allocated to the cluster that offers the earliest start-time.

task	tlevel	blevel
τ_4	0	7
τ_3	3	2
τ_1	0	5
τ_2	4	3

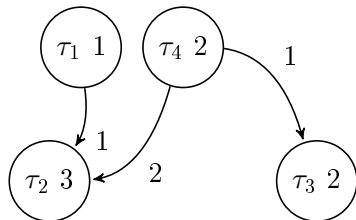


Figure: A Directed Acyclic Graph

DSC (Dominant Sequence Clustering)

[Yang and Gerasoulis 1994]

- $\text{priority}(\tau) = \text{tlevel}(\tau) + \text{blevel}(\tau)$.
- A zeroing(τ_p, τ) puts τ in the cluster of a predecessor $\tau_p \Rightarrow$ reduces $\text{tlevel}(\tau)$ by setting to zero the cost of the incident edge (τ_p, τ).

step	task	tlevel	blevel	DS	scheduled tlevel	
					κ_0	κ_1
1	τ_4	0	7	7	0*	
2	τ_3	3	2	5	2*	3
3	τ_1	0	5	5		0*
4	τ_2	4	3	7	5	4*

κ_0	κ_1
τ_4	τ_1
τ_3	τ_2

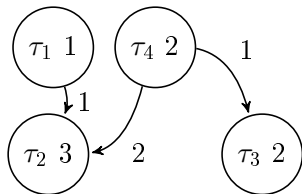


Figure: A Directed Acyclic Graph

- Number of processors is not predefined \rightarrow blind clustering.
- Memory size is not predefined \rightarrow blind clustering.
- Creates a new cluster when no zeroing is accepted \rightarrow creates long idle slots in already existing clusters.

\Rightarrow xDSC: A MEMORY-CONSTRAINED, NUMBER OF PROCESSOR-BOUNDED EXTENSION OF DSC

- ① Memory Constraint Warranty (MCW):
 - Verifying that the zeroing does not exceed a memory threshold M .
 - $\text{task_data}(\tau)$ is an overapproximation of the amount of memory used by Task τ .
 - $\text{cluster_data}(k)$ is an overapproximation of the amount of memory used by Cluster k .
 - $\text{size_data}(\text{data_merge}(\text{cluster_data}(k), \text{task_data}(\tau))) \leq M$.
- ② Bounded number of processors:
 - Verifying that new allocations do not exceed a cluster number threshold P .
 - $\text{cluster_time}(k)$ is the start time of the last scheduled task in k + its task_time .
 - otherwise, $\text{argmin}_{k \in \text{clusters}} \text{cluster_time}(k)$ under the constraint MCW.

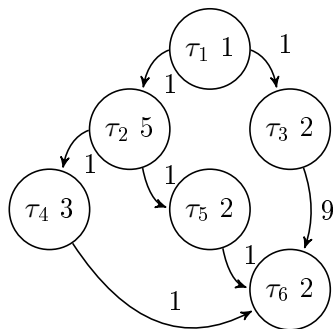


Figure: A DAG amenable to cluster minimization

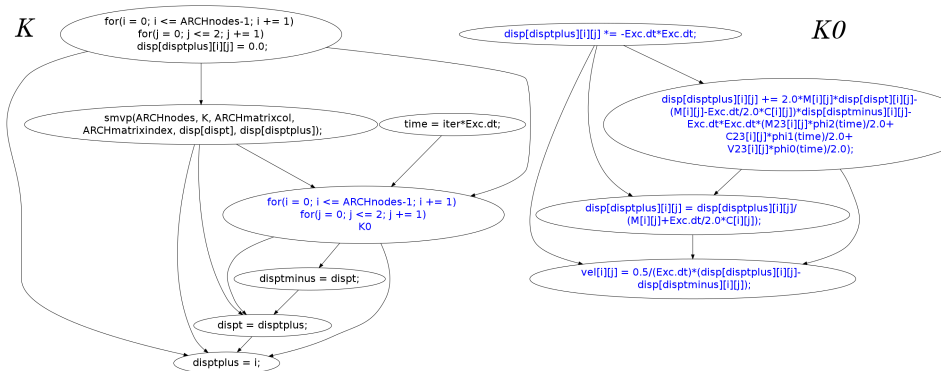
κ_0	κ_1	κ_2
τ_1		
τ_3	τ_2	
	τ_4	
τ_6		τ_5

- Allocation of τ to the last idle slot of κ ,
- Decreases $tlevel(\tau)$.
- For all nodes τ_s in κ :
 - $scheduled(successors(\tau_s))$,
 - $successors(\tau)$ are included in $successors(\tau_s)$.

step	task	t level	b level	DS	tlevel	
					κ_0	κ_1
1	τ_1	0	15	15	0*	
2	τ_3	2	13	15	1*	
3	τ_2	2	12	14	3	2*
4	τ_4	8	6	14		7*
5	τ_5	8	5	13	8*	10
6	τ_6	13	2	15	10*	

κ_0	κ_1
τ_1	
τ_3	τ_2
τ_5	τ_4
τ_6	

DAG: Hierarchical Clustered Dependence Graph (KDG)



- A test (both branches: true+false) constitutes one node (task).
 - A loop nest is an indivisible node.
 - A simple instruction is an indivisible node.
- ⇒ Hierarchy: recursively include KDGs.

Edge Cost, Task Time and Used Data

From Convex polyhedra to Polynomials

1 Edge Cost:

- Number of bytes of dependences RAW to annotate edges in the KDG,
- $\text{edge_cost}(\tau_i, \tau_j) = \text{size_data}(\text{regions_intersection}(\text{read_regions}(\tau_i), \text{write_regions}(\tau_j)))$.

2 Task Data:

- $\text{task_data}(\tau) = \text{data_merge}(\text{read_regions}(\tau), \text{write_regions}(\tau))$
- $\text{data_merge}(\mathbf{R}_1, \mathbf{R}_2) = \text{regions_union}(\mathbf{R}_1, \mathbf{R}_2) - \text{regions_intersection}(\mathbf{R}_1, \mathbf{R}_2)$

3 Size of regions (convex polyhedra) \Rightarrow Ehrhart polynomials represent the number of integer points contained in a given parameterized polyhedron.

4 Task Time:

- An estimation of complexity for each node in the KDG,
- $\text{task_time}(\tau) = \text{complexity_estimation}(\tau) \Rightarrow$ Polynomials.

① Applications

- Signal processing application ABF (Adaptive Beam Forming) [Griffiths 1969].
- Image processing application Harris corner detector [Harris and Stephens 1988]: detect the point of interest in an image.
- SPEC benchmark quake [Bao et al. 1998]: simulation of seismic wave propagation in large valleys.

② Machines

- SMP: 2-socket AMD quadcore Opteron with 8 cores, M = 16Gb of RAM, 2.4 GHz.
- DMP: 6 bicore processors Intel(R) Xeon(R), M = 32Gb of RAM per processor, 2.5 GHz.

From Polynomial to Values

Simple Cases

- When data are known numerical parameters, then each task polynomial is a constant (case of the application ABF).
- However, when input data are unknown at compile time (case of the application Harris), we use a simple heuristic to check the behavior of that polynomials, by comparing the coefficients of their monomials.
- Assume that all polynomials are monomials on the same bases.

Function	Complexity (polynomial)	Static time estimation
InitHarris	$9 \times \text{sizeN} \times \text{sizeM}$	9
SobelX	$60 \times \text{sizeN} \times \text{sizeM}$	60
SobelY	$60 \times \text{sizeN} \times \text{sizeM}$	60
Multiply	$20 \times \text{sizeN} \times \text{sizeM}$	20
Gauss	$85 \times \text{sizeN} \times \text{sizeM}$	85
CoarsitY	$34 \times \text{sizeN} \times \text{sizeM}$	34
One image transfer	$4 \times \text{sizeN} \times \text{sizeM}$	4

From Polynomial to Values

Instrumentation

```
FILE *finstrumented = fopen("instrumented_equake.in","w");
...
fprintf(finstrumented, "62_=%ld\n", 179 * ARCHElems + 3);
for (i = 0; i < ARCHElems; i++){
    for (j = 0; j < 4; j++){
        cor[j] = ARCHvertex[i][j];
    }
}
...
fprintf(finstrumented, "163_=%ld\n", 20 * ARCHnodes + 3);
for(i = 0; i <= ARCHnodes-1; i += 1)
    for(j = 0; j <= 2; j += 1)
        disp[disptplus][i][j] = 0.0;
fprintf(finstrumented, "163_>_166_=%ld\n", ARCHnodes * 9); //edge_cost(163,166)
fprintf(finstrumented, "166_=%ld\n", 110 * ARCHnodes + 106); //task_time(166)
smvp_opt(ARCHnodes, K, ARCHmatrixcol, ARCHmatrixindex, disp[dispt], disp[disptplus]);
fprintf(finstrumented, 167, 6);
time = iter*Exc.dt;
fprintf(instrumented_equake, 168, 510.50 * ARCHnodes + 3);
for (i = 0; i < ARCHnodes; i++){
    for (j = 0; j < 3; j++){
        disp[disptplus][i][j] *= -Exc.dt*Exc.dt;
        disp[disptplus][i][j] +=
            2.0*M[i][j]*disp[dispt][i][j]-M[i][j]-Exc.dt/2.0*C[i][j]*disp[disptminus][i][j] -
            Exc.dt * Exc.dt * (M23[i][j] * phi2(time) / 2.0 +
            C23[i][j] * phi1(time) / 2.0 + V23[i][j] * phi0(time) / 2.0);
        disp[disptplus][i][j] = disp[disptplus][i][j] / (M[i][j] + Exc.dt / 2.0 * C[i][j]);
        vel[i][j] = 0.5 / Exc.dt * (disp[disptplus][i][j] - disp[disptminus][i][j]);
    }
}
fprintf(finstrumented, "175_=%d\n", 2);
disptminus = dispt;
fprintf(finstrumented, "176_=%d\n", 2);
dispt = disptplus;
fprintf(finstrumented, "177_=%d\n", 2);
disptplus = i;
```

Experiments: ABF and equake

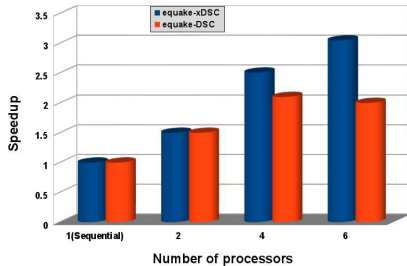
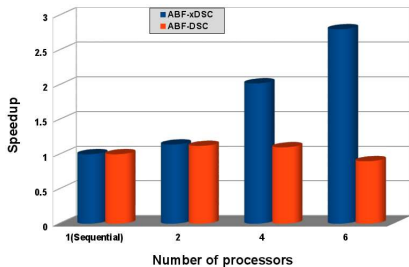
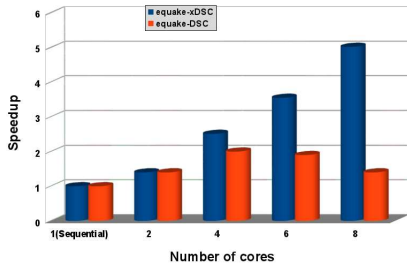
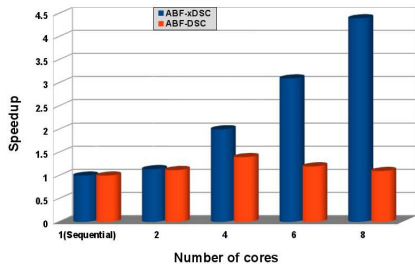
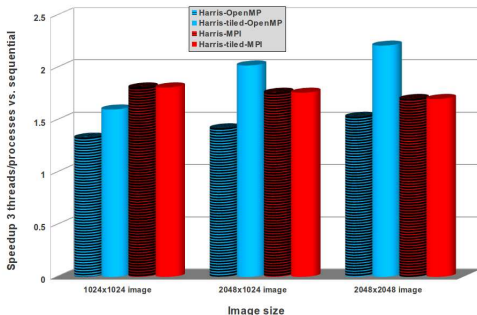
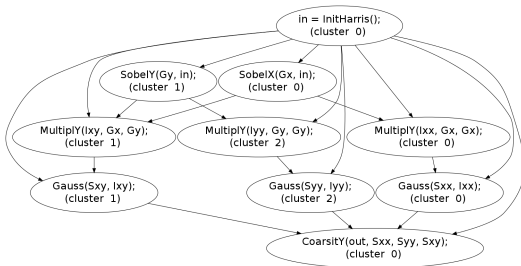


Figure: OpenMP/MPI vs. sequential speedup (ABF)

Figure: OpenMP/MPI vs. sequential speedup (equake)

Experiments: Harris



Conclusion

- xDSC: a new static scheduling,
- Precise and efficient cost model,
- Targeting both shared and distributed memory architectures,
- Memory constraint, Bounded number of processors, Efficient processor allocation.

Future Work

- Automatic code generation : OpenMP + MPI.
- Efficient hierarchical processor allocation strategy in order to yield a better xDSC-based parallelization process

Parallelizing with xDSC, a Resource-Constrained Scheduling Algorithm for Shared and Distributed Memory Systems

Dounia KHALDI

CRI, Mathématiques et systèmes
MINES ParisTech

June 19, 2012

