Domain-Specific Abstractions and Performance Portability & Data Access Complexity: Revisiting

the Red/Blue Pebble game

P. (Saday) Sadayappan The Ohio State University

Acknowledgements

Collaborators

Gerald Baumgartner(LSU) Albert Cohen (ENS Paris) Jason Cong (UCLA) Franz Franchetti (CMU) Robert Harrison (Stony Brook) So Hirata (U. Illinois) Jarek Nieploha (PNNL) Srini Parthasarathy (OSU) Louis-Noel Pouchet (UCLA) Russ Pitzer (OSU, Chem) J. Ramanujam (LSU) Fabrice Rastello (ENS Lyon) Nasko Rountev (OSU) Vivek Sarkar (Rice)

Ph.D. Students

Muthu Baskaran Uday Bondhugula Jim Dinan Xiaoyang Gao Albert Hartono Justin Holewinski Sriram Krishnamoorthy Qingda Lu Mohammad Arafat Venmugil Elango **Tom Henretty** Martin Kong Pai-Wei Lai Mahesh Ravishankar Kevin Stock Sanket Tavarageri

Funding

Natl. Science Foundn Dept. of Energy DARPA

Why Domain-Specific Languages?

- Productivity
 - High level abstractions eases application development
- Performance
 - Domain-specific semantics enables specialized optimizations
 - Constraints on specification enables more effective general-purpose transformations and tuning (tiling, fusion)
- Portability
 - New architectures => changes only in domain-specific compiler, without any change in user application code



DSL Technology for Exascale Computing (D-TEC)

Lead PI Daniel J. Quinlan Lawrence Livermore National Laboratory

Deputy PI Saman Amarasinghe MIT



The Tensor Contraction Engine A Domain-Specific Compiler for Many-Body Methods in Quantum Chemistry

Oak Ridge National Laboratory David E. Bernholdt, Robert Harrison

Pacific Northwest National Laboratory Jarek Nieplocha

Louisiana State University Gerald Baumgartner J. Ramanujam

Ohio State University

Xiaoyang Gao, Albert Hartono, Sriram Krishnamoorthy, Qingda Lu, Alex Sibiryakov, *Russell Pitzer, P. Sadayappan*

University of Florida So *Hirata*

University of Waterloo Marcel Nooijen

Supported by NSF and DOE

Time Crunch in Quantum Chemistry

Two major bottlenecks in computational chemistry:

- Very computationally intensive models
- Extremely time consuming to develop codes
- The vicious cycle of computational science:
- More powerful computers make more accurate models computationally feasible :-)
- But efficient parallel implementation of complex models takes longer and longer
- Hence computational scientists spend more time with low-level programming for performance, and less time doing science :-(
- Coupled Cluster family of models in electronic structure theory
- Increasing number of terms => explosive increase in code complexity
- Theory is the same, but efficient implementations of higher order models took many years

Theory	#Terms	#F77Lines	Year
CCD	11	3209	1978
CCSD	48	13213	1982
CCSDT	102	33932	1988
CCSDTQ	183	79901	1992

CCSD Doubles Equation (Quantum Chemist's Eye Test Chart :-)

 $hbar[a,b,i,j] == sum[f[b,c]*t[i,j,a,c], \{c\}] - sum[f[k,c]*t[k,b]*t[i,j,a,c], \{k,c\}] + sum[f[a,c]*t[i,j,c,b], \{c\}] - sum[f[k,c]*t[k,a]*t[i,j,c,b], \{c\}] - sum[$ $\{k,c\}$ -sum[f[k,j]*t[i,k,a,b], $\{k\}$ -sum[f[k,c]*t[i,c]*t[i,k,a,b], $\{k,c\}$ -sum[f[k,i]*t[i,k,b,a], $\{k\}$ -sum[f[k,c]*t[i,c]*t[i,k,b,a], $\{k\}$ -sum[f[k,c]*t[i,c]*t[$\{k,c\}$ + sum[t[i,c]*t[j,d]*v[a,b,c,d], {c,d}] + sum[t[i,j,c,d]*v[a,b,c,d], {c,d}] + sum[t[j,c]*v[a,b,i,c], {c}] - sum[t[k,b]*v[a,k,i,j], {c,d}] $\{k\}$ + sum[t[i,c]*v[b,a,j,c], {c}] - sum[t[k,a]*v[b,k,j,i], {k}] - sum[t[k,d]*t[i,j,c,b]*v[k,a,c,d], {k,c,d}] $sum[t[i,c]*t[j,k,b,d]*v[k,a,c,d], \{k,c,d\}] - sum[t[j,c]*t[k,b]*v[k,a,c,i], \{k,c\}] + 2*sum[t[j,k,b,c]*v[k,a,c,i], \{k,c\}] - 2*sum[t[j,k,b,c]*v[k,a,c], \{k,c\}] - 2*sum[t[j,k,b,c]*v[k,a,c], \{k,c\}] - 2*sum[t[j,k,b,c]*v[k,a,c], \{k,c\}] - 2*sum[t[j,k,b,c]*v[k,a,c], \{$ $sum[t[i,k,c,b]*v[k,a,c,i], \{k,c\}] - sum[t[i,c]*t[i,d]*t[k,b]*v[k,a,d,c], \{k,c,d\}] + 2*sum[t[k,d]*t[i,i,c,b]*v[k,a,d,c], \{k,c,d\}] - 2*sum[t[k,d]*t[k,d]*v[k,d]*v[k,d]*v[k,d]*v[k,d]*v[k,d]*v[k,d]*v[k,d]*v[k,d]*v[k,$ $sum[t[k,b]*t[i,j,c,d]*v[k,a,d,c], \{k,c,d\}] - sum[t[i,d]*t[i,k,c,b]*v[k,a,d,c], \{k,c,d\}] + 2*sum[t[i,c]*t[i,k,b,d]*v[k,a,d,c], \{k,c,d\}]$ $sum[t[i,c]*t[j,k,d,b]*v[k,a,d,c], \{k,c,d\}] - sum[t[j,k,b,c]*v[k,a,i,c], \{k,c\}] - sum[t[i,c]*t[k,b]*v[k,a,j,c], \{k,c\}] - sum[t[i,c]*t[k,b]*v[k,a,j], \{k,c\}] - sum[t[i,c]*t[k,b]*v[k,a,j], \{k,c]} - sum[t[i,c]*t[k,b]*v[k,a,j], \{k,c]\} - sum[t[i,c]*t[k,b]*v$ $sum[t[i,k,c,b]*v[k,a,j,c], \{k,c\}] - sum[t[i,c]*t[j,d]*t[k,a]*v[k,b,c,d], \{k,c,d\}] - sum[t[k,d]*t[i,j,a,c]*v[k,b,c,d], \{k,c,d\}] - sum[t[i,c]*t[j,d]*t[k,a]*v[k,b,c,d], \{k,c,d\}] - sum[t[k,d]*t[i,j,a,c]*v[k,b,c,d], \{k,c,d\}] - sum[t[k,d]*t[k,$ $sum[t[k,a]*t[i,j,c,d]*v[k,b,c,d], \{k,c,d\}] + 2*sum[t[i,d]*t[i,k,a,c]*v[k,b,c,d], \{k,c,d\}] - sum[t[i,d]*t[i,k,c,a]*v[k,b,c,d], \{k,c,d\}] - sum[t[i,d]*t[i,k,c,d], \{k,c,d\}] - sum[t[i,d]*t[i,k,c,d], \{k,c,d\}] - sum[t[i,d]*t[i,k,c,d], \{k,c,d], \{k,c,d\}] - sum[t[i,d]*t[i,k,c,d], \{k,c,d], \{k,c,d\}]$ $sum[t[i,c]*t[j,k,d,a]*v[k,b,c,d], \{k,c,d\}] - sum[t[i,c]*t[k,a]*v[k,b,c,j], \{k,c\}] + 2*sum[t[i,k,a,c]*v[k,b,c,j], \{k,c\}] - 2*sum[t[i,k,a,c]*v[k,b,c], \{k,c\}] - 2*sum[t[i,k,a,c]*v[k,b], \{k,c\}] - 2*sum[t[i,k,a,c], \{k,c\}] - 2*sum[t[i,k,a,c], \{k,c\}] - 2*sum[t[i,k,a,c], \{k,c$ $sum[t[i,k,c,a]*v[k,b,c,j],\{k,c\}] + 2*sum[t[k,d]*t[i,j,a,c]*v[k,b,d,c],\{k,c,d\}] - sum[t[j,d]*t[i,k,a,c]*v[k,b,d,c],\{k,c,d\}] - sum[t[j,d]*t$ $sum[t[j,c]*t[k,a]*v[k,b,i,c], \{k,c\}] - sum[t[j,k,c,a]*v[k,b,i,c], \{k,c\}] - sum[t[i,k,a,c]*v[k,b,j,c], \{k,c\}]$ $+sum[t[i,c]*t[j,d]*t[k,a]*t[l,b]*v[k,l,c,d], \{k,l,c,d\}] -2*sum[t[k,b]*t[l,d]*t[i,j,a,c]*v[k,l,c,d], \{k,l,c,d\}]$ $-2*sum[t[k,a]*t[1,d]*t[i,j,c,b]*v[k,1,c,d], \{k,1,c,d\}] + sum[t[k,a]*t[1,b]*t[i,j,c,d]*v[k,1,c,d], \{k,1,c,d\}]$ $-2*sum[t[j,c]*t[1,d]*t[i,k,a,b]*v[k,1,c,d], \{k,1,c,d\}] -2*sum[t[j,d]*t[1,b]*t[i,k,a,c]*v[k,1,c,d], \{k,1,c,d\}]$ $+sum[t[i,d]*t[i,k,c,a]*v[k,l,c,d],\{k,l,c,d\}] - 2*sum[t[i,c]*t[1,d]*t[i,k,b,a]*v[k,l,c,d],\{k,l,c,d\}]$ $+sum[t[i,c]*t[1,a]*t[j,k,b,d]*v[k,l,c,d], \{k,l,c,d\}] +sum[t[i,c]*t[1,b]*t[j,k,d,a]*v[k,l,c,d], \{k,l,c,d\}]$ $+sum[t[i,k,c,d]*t[j,l,b,a]*v[k,l,c,d], \{k,l,c,d\}] + 4*sum[t[i,k,a,c]*t[j,l,b,d]*v[k,l,c,d], \{k,l,c,d\}]$ $-2*sum[t[i,k,c,a]*t[i,l,b,d]*v[k,l,c,d], \{k,l,c,d\}] - 2*sum[t[i,k,a,b]*t[i,l,c,d]*v[k,l,c,d], \{k,l,c,d\}]$ $-2*sum[t[i,k,a,c]*t[i,l,d,b]*v[k,l,c,d], \{k,l,c,d\}] + sum[t[i,k,c,a]*t[i,l,d,b]*v[k,l,c,d], \{k,l,c,d\}]$ $+sum[t[i,c]*t[j,d]*t[k,l,a,b]*v[k,l,c,d], \{k,l,c,d\}] +sum[t[i,j,c,d]*t[k,l,a,b]*v[k,l,c,d], \{k,l,c,d\}]$ $-2*sum[t[i,j,c,b]*t[k,l,a,d]*v[k,l,c,d], \{k,l,c,d\}] -2*sum[t[i,j,a,c]*t[k,l,b,d]*v[k,l,c,d], \{k,l,c,d\}]$ $+sum[t[i,c]*t[k,b]*t[1,a]*v[k,1,c,i], \{k,1,c\}] +sum[t[1,c]*t[i,k,b,a]*v[k,1,c,i], \{k,1,c\}] -2*sum[t[1,a]*t[i,k,b,c]*v[k,1,c,i], \{k,1,c\}]$ $+sum[t[1,a]*t[j,k,c,b]*v[k,l,c,i], \{k,l,c\}] - 2*sum[t[k,c]*t[j,l,b,a]*v[k,l,c,i], \{k,l,c\}] +sum[t[k,a]*t[j,l,b,c]*v[k,l,c,i], \{k,l,c\}]$ $+sum[t[k,b]*t[j,l,c,a]*v[k,l,c,i], \{k,l,c\}] +sum[t[j,c]*t[l,k,a,b]*v[k,l,c,i], \{k,l,c\}] +sum[t[i,c]*t[k,a]*t[l,b]*v[k,l,c,i], \{k,l,c\}]$ $+sum[t[1,c]*t[i,k,a,b]*v[k,1,c,i], \{k,1,c\}] - 2*sum[t[1,b]*t[i,k,a,c]*v[k,1,c,i], \{k,1,c\}] +sum[t[1,b]*t[i,k,c,a]*v[k,1,c,i], \{k,1,c\}]$ $+sum[t[i,c]*t[k,l,a,b]*v[k,l,c,j], \{k,l,c\}] +sum[t[i,c]*t[l,d]*t[i,k,a,b]*v[k,l,d,c], \{k,l,c,d\}] +sum[t[i,d]*t[l,b]*t[i,k,a,c]*v[k,l,d,c], \{k,l,c,d\}]$ $\{k,l,c,d\}$ +sum[t[i,d]*t[1,a]*t[i,k,c,b]*v[k,l,d,c], $\{k,l,c,d\}$] -2*sum[t[i,k,c,d]*t[i,l,b,a]*v[k,l,d,c], $\{k,l,c,d\}$] $-2*sum[t[i,k,a,c]*t[i,l,b,d]*v[k,l,d,c], \{k,l,c,d\}] + sum[t[i,k,c,a]*t[i,l,b,d]*v[k,l,d,c], \{k,l,c,d\}]$ $+sum[t[i,k,a,b]*t[j,l,c,d]*v[k,l,d,c], \{k,l,c,d\}] +sum[t[i,k,c,b]*t[j,l,d,a]*v[k,l,d,c], \{k,l,c,d\}] +sum[t[i,k,a,c]*t[j,l,d,b]*v[k,l,d,c], \{k,l,c,d\}] +sum[t[i,k,a,c]*t[j,d,b]*v[k,d,c], \{k,l,c,d\}] +sum[t[i,k,a,c]*t[j,d,b]*v[k,d,c], \{k,l,c,d\}] +sum[t[i,k,a,c]*t[j,d,b]*v[k,d,c], \{k,l,c,d]\} +sum[t[i,k,c]*v[k,d,c], \{k,l,c], \{k,l,c], \{k,l,c,d\}] +su$ $\{k,l,c,d\} + sum[t[k,a]*t[l,b]*v[k,l,i,j], \{k,l\}] + sum[t[k,l,a,b]*v[k,l,i,j], \{k,l\}] + sum[t[k,b]*t[l,d]*t[i,j,a,c]*v[l,k,c,d], \{k,l,c,d\}]$ Γ/Γ1 ΠΦ/Γ1 ΠΦ/Γ¹ 1 ΠΦ Γ1 1 1 1 (1 1 1) Π Γ/Γ¹ ΠΦ/Γ1 ΠΦ/Γ¹ 1 ΠΦ/Γ¹ 1 1 (1 1 1)

Example: A CCSD Tensor Contraction

$$C_{ijkl} = \sum_{mn} A_{imkn} \cdot B_{jnlm}$$

for (i=0; i<P; i++)
for (j=0; j<Q; j++)
for (k=0; k<R; k++)
for (1=0; 1<S; 1++)
for (m=0; m<T; m++)
for (m=0; n<U; n++)
 C[i][j][k][1] += A[i][m][k][n]*B[j][n][1][m];</pre>

- Tensor contraction is a generalized high-dimension analog of matrix-matrix product
- Each loop index appears in exactly two tensors
 - Contraction indices appear only in the input (r.h.s) tensors: m and n
 - External indices appear in output tensor and one input tensor: {i, k} and {j, l}
- Each loop is either a parallel loop or a reduction loop
- Relative order of the contraction indices in loop nest is unimportant (associative reordering of reduction)
- Loop nests can be considered fully permutable (dependence analysis in existing compilers will not permit permutation of contraction indices)

Tensor Contraction Engine

- Automatic transformation from high-level specification
 - Chemist specifies computation in high-level mathematical form
 - Synthesis system transforms it to efficient parallel program
 - Code is tailored to target machine
 - Code can be optimized for specific molecules being modeled
- Multi-institutional collaboration (OSU, LSU, Waterloo, ORNL, PNNL, U. Florida)
- Two versions of TCE developed
 - a) Full chemistry, but fewer optimizations (Hirata) b) Excluded some details, but sophisticated optimizations
 - Used to implement over 20 models, in latest release of NWChem (a few million lines of synthesized code)
 - First parallel implementation for many of the methods
 - New improved TCE-2 planned

$A3A = \frac{1}{2} (X_{ce,af} Y_{ae,cf} + X_{c\bar{e},a\bar{f}} Y_{a\bar{e},c\bar{f}} + X_{c\bar{e},\bar{a}f} Y_{\bar{a}}$	e,cf
$+ X_{\overline{c}e,a\overline{f}}Y_{ae,\overline{c}\overline{f}} + X_{\overline{c}e,\overline{a}f}Y_{\overline{a}e,\overline{c}f} + X_{\overline{c}\overline{e},\overline{a}\overline{f}}Y_{\overline{a}\overline{e},\overline{c}}$	$\overline{f})$
$X_{ce,af} = t_{ij}^{ce} t_{ij}^{af} \qquad Y_{ae,cf} = \langle ab \ ek \rangle \langle cb \ fk$	$\left \right\rangle$
range V = 3000;	
range O = 100;	
index a,b,c,d,e,f : V;	
index i,j,k : O;	
mlimit = 1000000;	
function F1(V,V,V,O);	
function $F2(V,V,V,O)$;	
procedure F (in T1[O,O,V,V], in T2[O,O,V,V], out X)	=
begin	
$A3A == sum[sum[F1(a,b,e,k) * F2(c,f,b,k), {b,k}]$	
* sum[T1[i i c e] * T2[i i a f] {i i}]	
$\{a.e.c.f\}^{*}0.5 + \dots$	
end	

Customized Kernel Generation for Tensor Contractions

- Effective SIMD utilization is increasingly important for high performance on current/ emerging processors
- Automatic vectorization by production compilers (even with manual unrolling) often results in performance well under 50% of machine peak
- Customized code generator (using vector intrinsics) for tensor contractions

Approach to Code Generation

Overall Approach

- Explicitly generate vector-intrinsics based code for TC
- Compile generated code using icc/gcc
- Use Machine Learning model to predict performance of generated assembly code
- Explore space of code variants and choose the one with highest predicted performance

Search Space

Vectorized Dimension Determines how memory is accessed Loop Permutation Enable hoisting loads/stores outside inner loops Unroll-and-Jam Enable register reuse

Example: Multi-resolution Kernel

Kernel

$$R_{ijk} = \sum_{i'j'k'} S_{i'j'k'} X_{i'i} Y_{j'j} Z_{k'k}$$

- Used extensively in MADNESS (Multi-resolution Adaptive Numerical Environment for Scientific Simulation)
- Tensors are small, frequently fitting completely within L1-cache

Low Rank Decomposition $X_{i'i} = \sum_{l} X_{il}^{L} X_{li'}^{R}$ $Y_{j'j} = \sum_{m} Y_{jm}^{L} Y_{mj'}^{R}$ $Z_{k'k} = \sum_{n} Z_{kn}^{L} Z_{nk'}^{R}$

Example: Multi-resolution Kernel

Kernel

 $R_{ijk} = \sum_{i'i'k'} S_{i'j'k'} X_{i'i} Y_{j'j} Z_{k'k}$

 Implemented as a series of six tensor contractions

 $A_{ljk} = \sum_{i} S_{ijk} \cdot X_{il}^{L}$ $B_{lmk} = \sum_{i} A_{ljk} \cdot Y_{jm}^{L}$ $C_{lmn} = \sum_{k} B_{lmk} \cdot Z_{kn}^{L}$ $D_{lmk'} = \sum C_{lmn} \cdot Z^R_{nk'}$ $E_{lj'k'} = \sum D_{lmk'} \cdot Y^R_{mj'}$ $R_{i'j'k'} = \sum_{l} E_{lj'k'} \cdot X_{li'}^R$

Vectorization Dimensions







What if No Vectorizable Loop?

What if No Vectorizable Loop?

We can still vectorize by use of in-register transpose

Cost of register transpose often amortizable

- Vectors along columns are desired
- Vectors can only be loaded by rows



Code Example

1: p	rocedure IKJ(A _{ki} , B _{ik} , C _{ii})
2.	(j j.,
3:	for $(i \leftarrow 0; i < M; i + +)$ do
4:	for $(k \leftarrow 0; k < K; k + = 4)$ do
5:	$a_0[0:3] \leftarrow \text{SPLAT}(A[k+0][i])$
6:	$a_1[0:3] \leftarrow \operatorname{SPLAT}(A[k+1][i])$
7:	$a_2[0:3] \leftarrow \operatorname{SPLAT}(A[k+2][i])$
8:	$a_3[0:3] \leftarrow \operatorname{SPLAT}(A[k+3][i])$
9:	for $(j \leftarrow 0; j < N; j+=4)$ do
10:	$b_0[0:3] \leftarrow B[j+0][k:k+3]$
11:	$b_1[0:3] \leftarrow B[j+1][k:k+3]$
12:	$b_2[0:3] \leftarrow B[j+2][k:k+3]$
13:	$b_3[0:3] \leftarrow B[j+3][k:k+3]$
14:	TRANSPOSE (b_0, b_1, b_2, b_3)
15:	$c[0:3] \leftarrow C[i][j:j+3]$
16:	$c[0:3] + = a_0[0:3] * b_0[0:3]$
17:	$c[0:3] + = a_1[0:3] * b_1[0:3]$
18:	$c[0:3] + = a_2[0:3] * b_2[0:3]$
19:	$c[0:3] + = a_3[0:3] * b_3[0:3]$
20:	$C[i][j:j+3] \leftarrow c[0:3]$
21:	end for
55:	end for
54:	end for
<u> </u>	

Contraction $C_{ij} = \sum_{k} A_{ki} \cdot B_{jk}$

- Vectorized along j
- B_{jk} transposed
- Each element of A_{ki} is splatted (broadcast) to all elements of a vector register

Multiresolution Kernel Performance



Performance: CCSD Tensor Contraction

$$C_{ijkl} = \sum_{mn} A_{imkn} \cdot B_{jnlm}$$

Configuration	GCC	ICC	Machine Peak	Optimized
Nehalem double sse	1.406	2.740	10.64	7.579
Nehalem single sse	1.405	2.642	21.28	13.428
Sandy Bridge double avx	2.231	4.361	27.2	16.768
Sandy Bridge single avx	2.255	5.075	54.4	36.937
			1	

Performance in GFLOP/s

Domain-Specific Optimization: Stencils

Carnegie Mellon University Franz Franchetti, Richard Veras

Louisiana State University J. Ramanujam **Ohio State University**

Tom Henretty, Justin Holewinski, Martin Kong, *Nasko Rountev, P. Sadayappan*

UCLA Louis-Noel Pouchet

Supported by NSF and DOE

Embedded DSL for Stencils

- Benefits of high-level specification of computations
 - Ease of use
 - For mathematicians/scientists creating the code
 - Ease of optimization
 - Facilitate loop and data transformations by compiler
 - Automatic transformation by compiler into parallel C/C++ code
- Embedded DSL provides flexibility
 - Generality of standard programming language (C, MATLAB) for non compute-intensive parts
 - Automated transformation of embedded DSL code for high performance on different target architectures
- Target architectures for Stencil DSL

- Vector-SIMD (AVX, LRBNi, ..), GPU, FPGA, customized accelerators

Stencil DSL Example -- Standalone



Stencil DSL – Embedded in C

```
int main() {
 int Nr = 256; int Nc = 256; int T = 100;
 double *a = malloc(Nc*Nr*sizeof(double));
#pragma sdsl start time steps:T block:8,8,8 tile:1,3,1 time:4
 int Nr; int Nc;
 grid g [Nr][Nc];
 double griddata a on g at 0,1;
 pointfunction five point avg(p) {
    double ONE FIFTH = 0.2;
    [1]p[0][0] = ONE_FIFTH*([0]p[-1][0] + [0]p[0][-1])
               + [0]p[0][0] + [0]p[0][1] + [0]p[1][0]); }
 iterate 1000 {
    stencil jacobi 2d {
        [0:Nc-1] : [1]a[0][0] = [0]a[0][0];
      [0]
      [Nr-1][0:Nc-1]: [1]a[0][0] = [0]a[0][0];
      [0:Nr-1][0] : [1]a[0][0] = [0]a[0][0];
      [0:Nr-1][Nc-1]; [1]a[0][0] = [0]a[0][0];
      [1:Nr-2][1:Nc-2] : five point avg(a);}
   reduction max diff max {
      [0:Nr-1][0:Nr-1] : fabs([1]a[0][0] - [0]a[0][0]);
    }
  } check (max diff < .00001) every 4 iterations</pre>
 #pragma sdsl end
```

}

Related Work

- 20+ publications over the last few years on optimizing stencil computations
- Some stencil DSLs and stencil compilers
 - Pochoir (MIT), PATUS (Univ. Basel), Halide (MIT)
- Frameworks for building DSLs
 - SEJITS (LBL); Liszt, OptiML, OptiQL, … (Stanford)
- Our focus has been complementary: developing abstraction-specific compiler transformations matched to performance-critical characteristics of target architecture

Stencils on Vector-SIMD Processors

Fundamental source of

(e.g. SSE, AVX ...)

inefficiency with stencil codes on

current short-vector SIMD ISAs

Concurrent operations on

contiguous elements

```
for (i=0; i<H; ++i)
for (j=0; j<W; ++j)
c[i][j]+=b[i][j]+b[i][j+1];</pre>
```



Data in memory

Data Layout Transformation



- 1D vector in memory ⇔ (b) 2D logical view of same data
- (c) Transposed 2D array moves interacting elements into same slot of different vectors ⇔ (d) New 1D layout after transformation
- Boundaries need special handling

Data Layout Transformation: Evaluation



Standard Tiling with DLT





(b) Standard tiling -- DLT view (t=1)

- Standard tiling cannot be used with the layout transform
- Inter-tile dependences prevent vectorization

Split-tiling



Space

- Divide iteration space into upright and inverted tiles
- For each *tt* timesteps where *tt* = time tile size...
 - Execute all upright (U) tiles in parallel
 - Execute all inverted (I) tiles in parallel
- Nested split tiling can be used for multiple spatial dimensions
 - For 2D, 4 kinds of tiles: UU, UI, IU, II

Hybrid Split-tiling



- Trade-off between standard tiling (parallelogram shape) and split tiling (trapezoidal shape)
 - Split tiling enables inter-tile parallelism and DLT but cache footprint grows with time-tile size
 - Standard tiling inhibits inter-tile parallelism but cache footprint is independent of time-tile size
- Hybrid approach: combine split tiling along some inner spatial loops and standard tiling along outer spatial loops

Experimental Results: DLT+Split Tiling



Stencils on GPUs

- Vector-SIMD alignment problems non-existent
- Different optimization challenges: limited forms of synchronization, avoidance of thread divergence
- Overlapped tiling
 - Redundantly compute neighboring cells to avoid interthread-block sync, lower communication, and avoid thread divergence



Stencil Compiler for GPU: Performance



Multi-target Code Generation from SDSL



Many Stand-Alone DSL Efforts



- Our DSL compilers for tensors and stencils did not reuse components
- Brute-force approach: N*M compilers for N domains and M platforms
- Can we achieve reuse in optimizing many DSLs for multiple targets?

A Layered Approach

Stencil DSL		Ten D:	isor SL						DSL-k
Domain-Specific Transformations Architecture-Independent Domain-Specific Layer						Specific oners c Layer			
Polyhedral Gener Transformations Trans			eneral ransfo	ral-purpose General formations Partit			purpose oners		
Architecture-Independent General-Transformations Architecture-Specific Kernel Optimizers									
Multi-cor CPU	e		GI	թՈ		FPGA			
Codelet-Centric Vector-SIMD Optimization

- Problem: Compilers like icc typically achieve only a fraction of processor's peak performance on polyhedrally transformed code
 - Overheads due to unaligned load/store operations
 - Overheads due to repeated load/stores of reused data elements
- Solution: Develop a decoupled two-step optimization process that allows separation of concerns:
 - Use polyhedral compiler transformations to form small (L1 cache resident) tiles with specific desirable properties: vector codelet
 - Use a low-level code optimizer (SPIRAL back-end) to generate optimized code for vector codelet
- Properties of vectorizable codelet: formulated as an Int. Lin. Prog.
 - Maximal inner-most loop parallelism
 - Maximal number of stride 0/1 accesses in inner-most loop
 - Maximize number of innermost permutable loops
 - Minimization of unaligned load/store operations
- Same approach now being targeted customized accelerators
- Details in PLDI '13 paper

Experimental Results: Vector Codelets

covariance



IO Complexity: Revisiting the Red/Blue Pebble Game

ENS Fabrice Rastello Ohio State University Venmugil Elango, *P. Sadayappan*

Louisiana State University J. Ramanujam

UCLA Louis-Noel Pouchet

High-Level Summary

- Characterizing "communication" complexity is much more complicated than computational complexity
 - Number of data transfers from large "slow" memory to limited "fast" memory (cache, registers etc.)
 - Unlike comp. complexity, depends on order of execution of operations and amount of fast memory
- First IO lower bounds: Hong/Kung classic 1981 paper
 - Tight lower bounds for some algorithms but loose for others
 - Monolithic: Cannot compose analyses: S = S1;S2
 - Only used for analysis of regular Computational DAGs
 - Recent advance by Demmel, Yelick, and colleagues at Berkeley provides a generalization for a class of loop computations (linear and tensor algebra, N-body, etc.), but not effective for stencils etc
- We develop pebbling variant and new bounding technique
 - 1. Enables composability: LB(S) from LB(S1) and LB(S2), S = S1;S2
 - 2. Tighter lower bounds for complementary class of CDAGs
 - 3. Enables full automation for arbitrary, irregular CDAGs

CDAG characterization [Bilardi 2001]

A computation directed acyclic graph (CDAG) is a 4-tuple C = (I, V, E, O) of finite sets such that:

- 1. $I \cap V = \emptyset;$
- 2. $E \subseteq (I \cup V) \times V$ is the set of arcs;
- 3. $G = (I \cup V, E)$ is a directed acyclic graph with no isolated vertices;
- 4. *I* is called the input set;
- 5. V is called the operation set and all its vertices have one or two incoming arcs;
- 6. $O \subseteq (I \cup V)$ is called the output set.

Hong-Kung (1981): Red-Blue Pebble Game

Given a CDAG G, S red pebbles and an arbitrary number of blue pebbles, with a blue pebble on each *input* vertex, a complete game is any sequence of steps obeying the following rules that results a final state with blue pebbles on all *output* vertices:

- **R1 (Input)** A red pebble may be placed on any vertex that has a blue pebble (load from slow to fast memory),
- **R2 (Output)** A blue pebble may be placed on any vertex that has a red pebble (store from fast to slow memory),
- **R3 (Compute)** If all immediate predecessors of a vertex have red pebbles, a red pebble may be placed on that vertex (execution or "firing" of operation),
- R4 (Delete) A red pebble may be removed from any vertex (reuse storage).

Code Example and its CDAG




























































Hong-Kung (1981): S-partitioning of CDAG

Given a CDAG G. An S-partitioning of G is a collection of h subsets of \mathcal{V} such that:

P1
$$\bigcap_{i=1}^{h} \mathcal{V}_i = \emptyset$$
, and $\bigcup_{i=1}^{h} \mathcal{V}_i = \mathcal{V}$

 ${\bf P2}\,$ there is no circuit between subsets

P3 $\forall i, \exists D \in \mathsf{Dom}(\mathcal{V}_i) \text{ such that } |D| \leq S$

```
P4 \forall i, |\mathsf{Min}(\mathcal{V}f_i)| \leq S
```

Where

- a dominator set of \mathcal{V}_i , $D \in \mathsf{Dom}(\mathcal{V}_i)$ is a set of vertices such that any path from I to a vertex in \mathcal{V}_i contains some vertex in D;
- the minimum set of \mathcal{V}_i , $Min(\mathcal{V}_i)$ is the set of vertices in \mathcal{V}_i that have at least one child outside of \mathcal{V}_i ; and |Set| is the cardinality of the set Set.

Hong-Kung (1981): S-partitioning theorem

Any complete calculation of the red-blue pebble game on on a CDAG using at most S red pebbles, is associated with a 2S-partition of CDAG such that

 $S \times h \ge q \ge S \times (h-1)$

where q is the number of I/O moves in the game and h is the number of subsets in the 2S-partition.

Lower Bound

Let H(2S) be the minimal number of vertex sets for any valid 2S-partition of a given CDAG. Then the minimal number Q of I/O operations for any valid execution of the CDAG is bounded by:

$$Q \ge S \times (H(2S) - 1)$$

The Composability Problem



- The IO for a computation includes all inputs and outputs of the CDAG
- Consider S = S1 ; S2
- If some outputs of S1 are inputs to S2, they may be passed in red pebbles in a complete game for S
- OptIO(S) <= OPtIO(S1)+OPtIO(S2)
- Hence Lbound(S1) and Lbound(S2) cannot simply be added to get Lbound(S)
- In contrast Comp(S) = Comp(S1)+Comp(S2)

Addressing Composability

- 1. Modify pebble game to disallow re-computation
 - Add a third type of pebble: white
 - Placed white pebble on a vertex, along with red pebble, when vertex is first computed
 - Compute rule changed to only allow it when vertex does not have a white pebble (no change to load, store and delete rules)
- 2. Define a new metric IIO (internal IO) of CDAG G:
 - Given G = (I, V, E, O), internal CDAG G' = (O, V, E, O)
 - IIO(G) is IO(G'), for optimal game on G'
 - $IIO(G) \le IO(G) \le IIO(G) + |I| + |O|$

Composable
Lower Bounds $IIO(G1)+IIO(G2) \le IIO(G1 \cup G2)$ $IIO(G1)+IIO(G2) \le IO(G1 \cup G2)$ where as
IO(G1 U G2) \le IO(G1)+IO(G2) $IO(G1 \cup G2) \le IO(G1)+IO(G2)$ Lower Bounds
NOT composable

Red-Blue-White Pebble Game

Given a CDAG G = (I, V, E, O), S red pebbles and an arbitrary number of blue and white pebbles, with a blue and white pebble on each *input* vertex, a complete game is any sequence of steps obeying the following rules that results a final state with blue pebbles on all *output* vertices:

R1 (Input) A red pebble may be placed on any vertex that has a blue pebble,

R2 (Output) A blue pebble may be placed on any vertex that has a red pebble,

R3 (Compute) If a vertex v does not have a white pebble and all its immediate predecessors have red pebbles on them, a red pebble along with a white pebble may be placed on v.

R4 (Delete) A red pebble may be removed from any vertex (reuse storage).

S-partitioning of CDAG – new definition

Given a CDAG G. An S-partitioning of G is a collection of h subsets of $\mathcal{V} - I$ such that:

P1
$$\bigcap_{i=1}^{h} V_i = \emptyset$$
, and $\bigcup_{i=1}^{h} V_i = \mathcal{V} - I$

 $\mathbf{P2}$ there is no circuit between subsets

P3 $\forall i$, $|\ln(V_i)| \leq S$

P4 $\forall i$, $|\mathsf{Out}(V_i)| \leq S$

Where the input set of V_i , $\ln(V_i)$ is the set of vertices of $\mathcal{V} - V_i$ that have at least one child in V_i ; the output set of V_i , $Out(V_i)$ is the set of vertices of V_i that have at least one child outside of \mathcal{V}_i ; and |Set| is the cardinality of the set *Set*.





























Limitation of Hong-Kung S-partition Model

- Hong/Kung's S-partitioning approach constrains vertex partitions based on size of input dominator sets
 - Not effective for computations with few inputs that generate many more intermediates and finally output few values; extreme example: Diamond DAG
 - For S>0 red pebbles, a valid 2Spartition is the entire CDAG,
 i.e. lower bound = S*(1-1) = 0



New Approach: MinCut Partitioning

- Association between wavefront of "live" vertices just before firing vertex x in any valid RBW pebble game and a convex graph mincut including vertex x
- Number of vertices in maximal convex graphmincut associated with all vertices is a lower bound on the maximal *schedule* wavefront in optimal pebble game.



Diamond DAG with const. mincut wrt x



Constrained mincut wrt x . н. S



Diamond DAG with 4 disjoint vertex sets



IO Complexity Comparison

	Hong-Kung	Our bounds
Diamond DAG	0	$\frac{(n-2S)^2}{S}$
FFT	$\frac{\Omega(n\log(n))}{\log(S)}$	$\frac{2n\log(n)}{\log(S)}$
9-pt Seidel	$\Omega(\frac{n^2}{S})$	$0.75 \frac{n^2}{\sqrt{S}}$

IO Lower Bounds for Arbitrary CDAGs

- Automation: Tool to generate IO lower bound for an arbitrary CDAG instance, for given number of red pebbles
- Combine both approaches and use tighter bound:
 - 2S partitioning of Hong-Kung
 - New convex mincut
- Each bounding approach is effective for different class of CDAGs
 - 2S partitioning: high-fanout DAGs like Mat-Mult
 - Convex mincut: bounded fanout CDAGs like stencils

```
Function CalculateLB(G, do_mincut, do_2Kpart)
     if do mincut then
         Q_{mincut} = MinCut(G);
         if Q_{mincut} == 0 then
              do_mincut = false;
         end
     else
         Q_{mincut} = 0;
     end
     if do 2KPart then
         Q_{2KPart} = 2KPart(G);
         if Q_{2KPart} == 0 then
              do_2Kpart = false;
         end
     else
          Q_{2KPart} = 0;
    end
     if do_mincut v do_2Kpart then
         (G_1, G_2) = \text{Bisect}(G);
         Q_1 = CalculateLB(G_1, do_mincut, do_2Kpart);
         Q_2 = CalculateLB(G_2, do_mincut, do_2Kpart);
         return max (Q_{mincut}, Q_{2KPart}, Q_1, Q_2);
     else
         return 0;
     end
end
```

Summary

- Characterizing lower bounds on inherent data access complexity of algorithms is important
- Previous approaches have some limitations
- New approach to deriving lower bounds
 - Develop tighter analytical bounds for some algorithms
 - Enable composing lower bounds for constituent subcomputations
 - Generate lower bounds for irregular CDAGs