

The Data-Dependence graph of Adjoint Codes

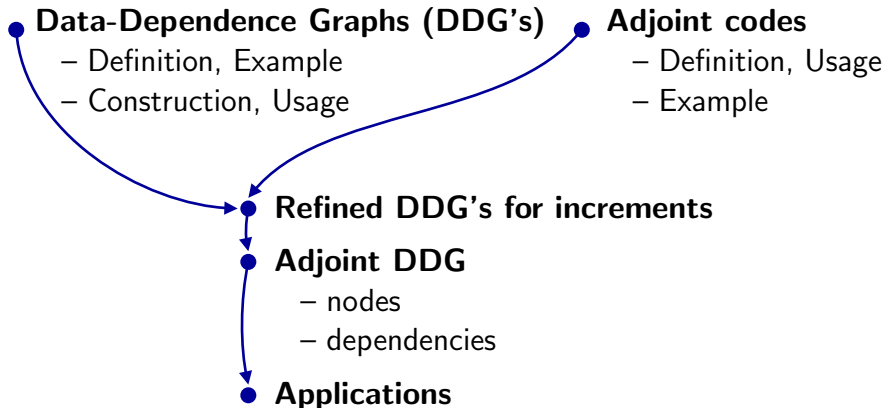
Laurent Hascoët

INRIA Sophia-Antipolis

November 19th, 2012

Data-Dependence graph of this talk

Relate the parallel properties of a code to those of its adjoint



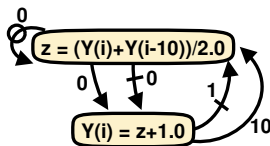
Data-Dependence Graphs (DDG's)

Given the source of a program,

- The DDG nodes are the operations of the program.
Granularity can be chosen: ops, instructions, blocks ...
- Iterated DDG nodes are collapsed
 \Rightarrow notion of distance in iteration spaces
- The DDG arrows are the necessary execution partial order.
Collapsed nodes \Rightarrow possible cycles!

```
X(:) = 0.0  
DO i=10,900  
  z = (Y(i)+Y(i-10))/2.0  
  Y(i) = z+1.0  
ENDDO
```

$X(:) = 0.0$



How DDG's are built

Data-Dependencies originate from variables (*"data"*).

A variable v causes a Data-Dependency from node $N1$ to node $N2$ iff

- $N1$ writes v and $N2$ reads v (*"true dependency"*), or
- $N1$ reads v and $N2$ overwrites v (*"anti dependency"*), or
- $N1$ writes v and $N2$ overwrites v (*"output dependency"*).

and in addition (*but not always necessary...*) v is not totally overwritten between $N1$ and $N2$.

What DDG's are for

A Data-Dependency from $N1$ to $N2$ means that any rescheduling of operations must keep the order $N1$ before $N2$.

⇒ DDG's are central for parallelization.

- DDG's capture notions of dead code, inlining ...
- In loop parallelization, DDG cycles (SCC) cannot be parallelized nor vectorized.
- Loop nests can be transformed according to Data-Dependency distances.
- In Message-Passing, DDG cycles capture deadlocks.
- Scalar expansion is a classical way to lift anti and output dependencies, thus breaking cycles.

Automatic Differentiation's Adjoint codes

- Adjoints are the most efficient way to obtain analytical gradients.
- Uses abound: optimization, least-squares, parameter estimation, sensitivities, uncertainty quantification.
- Adjoints can be built mechanically through AD.

Adjoints propagate backwards the gradient of the result. Consider P

$$P : \{l_1; l_2; \dots l_{p-1}; l_p\}$$

that computes a function $F : X \mapsto Y$

$$F(X) = Y = V_p = f_p(f_{p-1}(\dots f_2(f_1(V_0 = X)) \dots))$$

its gradient is

$$\frac{\partial Y}{\partial X} = f'_p(V_{p-1}) \times f'_{p-1}(V_{p-2}) \times \dots \times f'_2(V_1) \times f'_1(V_0)$$

that the adjoint code evaluates from left (row vector) to right.

Example

$$r = x * x + y * y$$

$$r = \text{sqrt}(r)$$

$$x = x / r$$

$$y = y / r$$

$$\begin{aligned}\bar{r} &= \bar{r} - \bar{y} * y / (r * r) \\ \bar{y} &= \bar{y} / r\end{aligned}$$

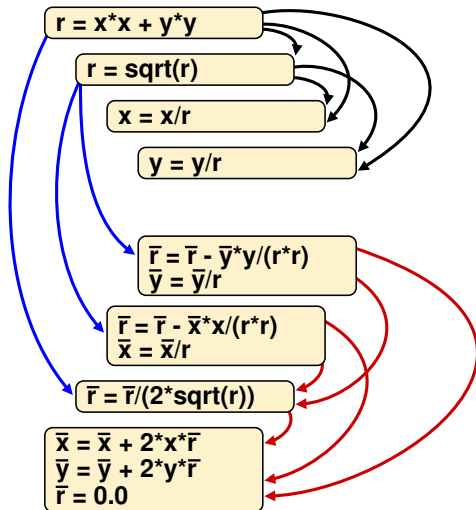
$$\begin{aligned}\bar{r} &= \bar{r} - \bar{x} * x / (r * r) \\ \bar{x} &= \bar{x} / r\end{aligned}$$

$$\bar{r} = \bar{r} / (2 * \text{sqrt}(r))$$

$$\begin{aligned}\bar{x} &= \bar{x} + 2 * x * \bar{r} \\ \bar{y} &= \bar{y} + 2 * y * \bar{r} \\ \bar{r} &= 0.0\end{aligned}$$

- Two successive parts (sweeps)
- 2nd sweep computes the gradient, reverse order.
- 1st sweep computes needed original values \Rightarrow copy of original code.
- Mechanism (not shown here) recovers needed r , x , y
- Gradient obtained in final \bar{x} and \bar{y} .

What about the DDG ?



- Chosen granularity: l_k and \bar{l}_k
- Black deps are the original deps on the original copy.
- Blue deps show uses of direct values for the derivatives.
- Red deps caused only by derivative variables.
- Red deps seem symmetric of Black deps !?

→ true ? why ?

Two remarks

1: For DDG, **increments** behave just like **reads**

2: The adjoint of a **read** is an **increment**, and conversely.

1: Increments behave like reads

- Define the effect of DDG node on variable v to be \textcircled{i} iff v is **only** used like $v = v + \dots$
- Warning: make sure increments are atomic !

Successive increments are **data-independent**
 \Rightarrow refined cases for Data-Dependency:

	\textcircled{n}	\textcircled{r}	\textcircled{i}	\textcircled{w}
\textcircled{n}				
\textcircled{r}			•	•
\textcircled{i}		•		•
\textcircled{w}		•	•	•

2: Adjoint of read is increment, and vice-versa

- Suppose l_k (only) reads v $(\textcircled{\mathbf{r}})$

- then $f'_k = \begin{pmatrix} \ddots & \bullet & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \ddots \end{pmatrix}$

- so that \overline{l}_k actually executes:

$$\begin{pmatrix} \dots & \overline{v} & \dots \end{pmatrix} = \begin{pmatrix} \dots & \overline{v} & \dots \end{pmatrix} \times \begin{pmatrix} \ddots & \bullet & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \ddots \end{pmatrix}$$

- which implies \overline{l}_k only increments \overline{v} $(\textcircled{\mathbf{i}})$
... or just doesn't mention it $(\textcircled{\mathbf{n}})$

Table of adjoint DDG nodes

- Looking at all cases:

- $N(\mathbf{n})$ on $v \Rightarrow \bar{N}\{\mathbf{n}\}$ on \bar{v}
- $N(\mathbf{r})$ on $v \Rightarrow \bar{N}\{\mathbf{i}, \mathbf{n}\}$ on \bar{v}
- $N(\mathbf{i})$ on $v \Rightarrow \bar{N}\{\mathbf{r}, \mathbf{n}\}$ on \bar{v}
- $N(\mathbf{w})$ on $v \Rightarrow \bar{N}\{\mathbf{w}, \mathbf{r}, \mathbf{i}, \mathbf{n}\}$ on \bar{v}

- which implies conversely:

- $\bar{N}(\mathbf{n})$ on $\bar{v} \Rightarrow N\{\mathbf{w}, \mathbf{r}, \mathbf{i}, \mathbf{n}\}$ on v
- $\bar{N}(\mathbf{r})$ on $\bar{v} \Rightarrow N\{\mathbf{w}, \mathbf{i}\}$ on v
- $\bar{N}(\mathbf{i})$ on $\bar{v} \Rightarrow N\{\mathbf{w}, \mathbf{r}\}$ on v
- $\bar{N}(\mathbf{w})$ on $\bar{v} \Rightarrow N\{\mathbf{w}\}$ on v

... and finally

If $\overline{l_k} \xrightarrow{\quad} \overline{l_j}$
then $l_j \xrightarrow{\quad} l_k$

- In other words, the adjoint DDG is equal to (or smaller than) the reversed original DDG
- Distances are preserved
- Nature of dependences (*true*, *anti*, *output*) is not preserved

Conclusion, Applications

- The adjoint DDG has the same structure as the original DDG
- Most parallel properties can propagate to the adjoint code:
 - The adjoint of a vector code is a vector code (in most cases).
 - The adjoint of a loop with independent iterations has independent iterations.
- There can even be **more** parallelism in the adjoint !
- ... but keep in mind the hypotheses (atomic increments).

Each HPC paradigm (e.g. Message-Passing SPMD) can use this DDG property in a specific way