

# Reasoning About Sparse Vectors for Loops Code Generation

Vadim Zaliva, Franz Franchetti, Carnegie Mellon University

## Motivating example

Pointwise operator vectorization:

$$f \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{pmatrix} = \begin{pmatrix} f(a_0) \\ 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ f(a_1) \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ f(a_2) \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \\ f(a_3) \end{pmatrix} = \begin{pmatrix} f(a_0) \\ f(a_1) \\ 0 \\ f(a_3) \end{pmatrix}$$

Loop Parallelization:

$$P_{f_j}^n = \sum_{j=0}^1 S_{(j)} \circ A_{f_j} \circ G_{(j)}$$

$$P_{f_j}^n = \sum_{j=0}^1 S_{h_{j,0}} \circ \left( \prod_{k=0}^1 A_{f_k} \right) \circ G_{h_{j,2}}$$

Sample non-vectorized function processing single FP element.

```
void f1(float *src, float *dst);
```

Pointwise application of 'f' to src, storing results in dst, one at a time. It requires 4 iterations and 4 function calls.

```
for(i=0;i<4;i++)
  f1(src+i,dst+i);
```

Sample vectorized function processing 2 FP elements at a time.

```
void f2(float *src, float *dst);
```

Pointwise application of 'f' to src, storing results in dst, two at a time. It requires just 2 iterations and 2 function calls.

```
for(i=0;i<2;i++)
  f2(src+2*i,dst+2*i);
```

## The Problem

- Dense vectors are decomposed into iterative sums of sparse vectors.
- Various decompositions (number or vectors and location of non-sparse values) could represent a variety of memory access patterns.
- This allows applying a variety of algebraic transformations to reshape a computation to optimize for vectorization, parallelization, sequential memory access.
- However in such iterative sum, the addition has a special semantics:
  - Mathematically, the sparse values could be treated as zeroes.
  - Operationally, combining sparse and non-sparse value could be seen as an assignment.
- The expressions produced are naturally mapped to SSA form only if certain constraints on structure of sparse vectors under iterative sums are maintained.
- Tracking and enforcing such constraints for correctness proofs is difficult, as they are not adequately enforced by mathematical abstraction used.
- In this work we present a working approach for structural constraints tracking and propagation used in Coq proofs of correctness.

## Approach

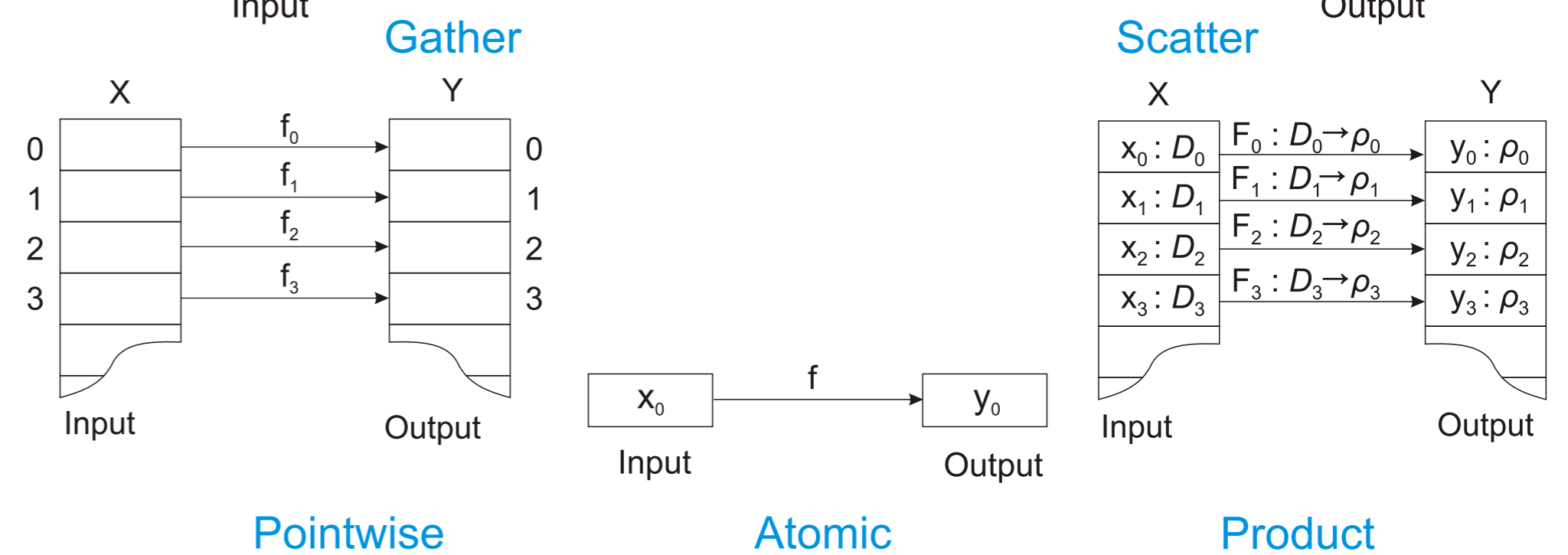
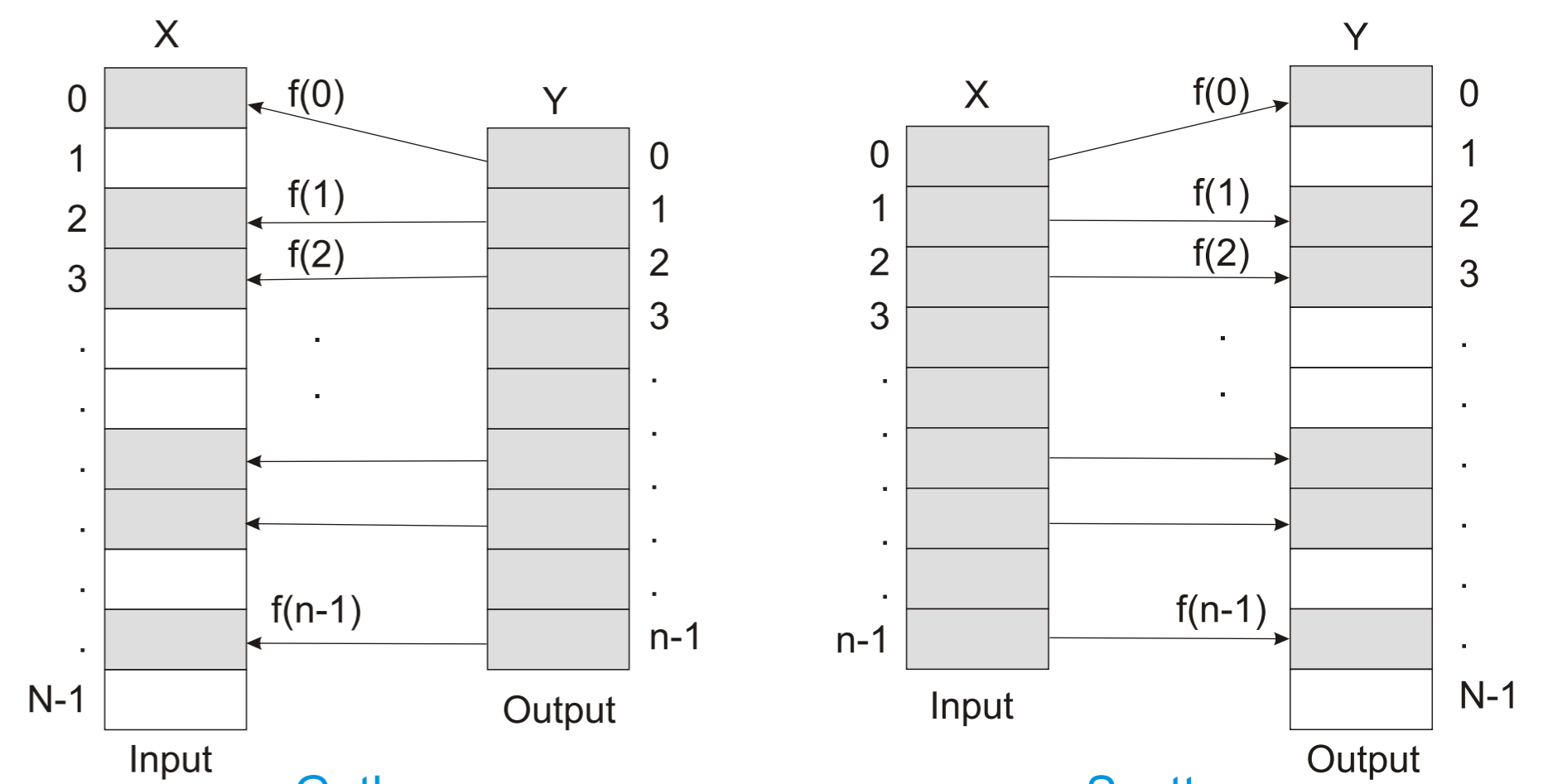
### Sparsity Requirements

- Distinguish empty and assigned cells.
- Treat empty cells as some "default" value. Such default value could depend on the context (e.g. 0 for addition but 1 for multiplication).
- In case of "a loop as sparse vector sum" we should never attempt to combine two non-sparse elements. This type of error we will call a "collision".
- Sparsity tracking should be easy to perform during computations.
- Collisions should be seamlessly tracked and propagated across the computation.
- The collision and sparsity tracking should be proof-friendly (easy to deal with in Coq)
- Separate sparsity tracking from actual operations on values as they represent two different aspects of computation.

### State and Collision Tracking Monad

```
Record Rflags : Type := mkRFlags {is_struct: bool ; is_collision: bool }.
Definition rFlagsZero := mkRFlags true false.
Definition mappend (a b: Rflags): Rflags :=
  mkRFlags
    (is_struct a && is_struct b)
    (is_collision a || (is_collision b ||
      (negb (is_struct a || is_struct b)))).
Definition RMonoid : Monoid Rflags := Build_Monoid mappend rFlagsZero.
Definition R_theta := writer Rmonoid R.
```

## Operators



## Implementation

### "Diamond" Abstraction

Scalar  $\diamond : \mathcal{A} \rightarrow \mathcal{A} \rightarrow \mathcal{A}$

Vector  $\vec{\diamond} : \mathcal{A}^n \rightarrow \mathcal{A}^n \rightarrow \mathcal{A}^n$   
 $((a_0, a_1, \dots, a_{n-1}), (b_0, b_1, \dots, b_{n-1})) \mapsto (a_0 \diamond b_0, a_1 \diamond b_1, \dots, a_{n-1} \diamond b_{n-1})$

Operation  $\delta : (\mathcal{A}^n \rightarrow \mathcal{A}^m) \rightarrow (\mathcal{A}^n \rightarrow \mathcal{A}^m) \rightarrow (\mathcal{A}^n \rightarrow \mathcal{A}^m)$   
 $(F, G) \mapsto (x \mapsto F(x) \delta G(x))$

Iterative  $\diamond_{i=0}^{n-1} F_i : \mathcal{A}^n \rightarrow \mathcal{A}^n$   
 $x \mapsto \begin{cases} \mathbf{0}^n & \text{if } n = 0, \\ (F_{n-1} \delta (\diamond_{j=0}^{n-2} F_j))(x) & \text{otherwise} \end{cases}$

### Iterative Sum with Sparsity and Collision Tracking

Let us apply the diamond abstraction demonstrated to  $R_\theta$  type (which represents  $\mathbb{R}$  with  $R_{flags}$  state) and summation operator. To do so we specialize previous notation as follows:

```
Definition A := R_theta.
Definition diamond := liftM2 (+).
Definition vec_diamond := vector.map2 diamond.
Definition delta f g := lambda x => (f x) vec_diamond (g x).
Definition O^n := vector.const (ret 0) n.
```

This gives us a sparse, collision-tracking Pointwise:

$$\text{Pointwise}_{n,f} = \diamond_{j=0}^{n-1} (S_{(j)} \circ A_{f_j} \circ G_{(j)})$$

## Summary

- Implementation in Coq proof assistant.
- Each value is tagged with two boolean flags: is\_struct and is\_collision.
- Flags structure along with combining operation forms a Monoid.
- Two Monoid instances are used: with and without collision tracking.
- Flags are tracked using Writer Monad.
- Operations on values can not directly examine sparsity flags and thus can not depend on them.
- Sparsity is automatically tracked by the monad. No implicit flags handling in operators implementation.
- Collision is automatically tracked and propagated by the monad.

## Contact info

Vadim Zaliva <vzaliva@cmu.edu>  
 SPIRAL: <http://spiral.net/>

Carnegie Mellon

