



Lightweight Polytypic Staging: a new approach to Nested Data Parallelism in Scala



Alexander Slesarenko
Keldysh Institute of Applied Mathematics, 2012

The Domain - Nested Data Parallelism

- ▶ The original idea
 - ▶ Guy Blelloch, Gary Sabot: in the early 90's ([1] is a good starting point)
 - ▶ NESL – proof of concept (first order, interpreted language)
- ▶ Generalizations (90's – 00's)
 - ▶ Chakravarty, Keller, S. P. Jones et al. (5 or 6 papers)
 - ▶ Data Parallel Haskell – higher-order, compiled language [2]
 - ▶ Language extension with special syntax
- ▶ A big promise but still in research

[1] Guy E. Blelloch. *Vector models for data-parallel computing*. MIT Press, Cambridge, MA, USA, 1990

[2] Simon Peyton Jones, Roman Leshchinskiy, Gabriele Keller, and Manuel M. T. Chakravarty. *Harnessing the Multicores: Nested Data Parallelism in Haskell*, 2008

Motivation: NDP as an embedded DSL

- ▶ NDP is not a “silver bullet”
- ▶ Some applications **fit** to the model but others **don't**
- ▶ For those that **fit** we want high-level declarative language
- ▶ **IDEALLY: If it is expressible then it is automatically vectorizable** (with asymptotic work-efficiency)
- ▶ Should interact with other DSLs and the host language
- ▶ Yet another tool in the Scala toolbox

Polytypic DSL

DSL	Monotypic (traditional)	Polytypic (data type generic)
Shallow embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Execution on the JVM	
Deep embedding		

Polytypic DSL

DSL	Monotypic (traditional)	Polytypic (data type generic)
Shallow embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Execution on the JVM	
Deep embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Staging + transform✓ Execution on XXX by code generation	

Polytypic DSL

DSL	Monotypic (traditional)	Polytypic (data type generic)
Shallow embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Execution on the JVM	<ul style="list-style-type: none">✓ Type-indexed types and functions✓ Execution on the JVM
Deep embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Staging + transform✓ Execution on XXX by code generation	

Polytypic DSL

DSL	Monotypic (traditional)	Polytypic (data type generic)
Shallow embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Execution on the JVM	<ul style="list-style-type: none">✓ Type-indexed types and functions✓ Execution on the JVM
Deep embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Staging + transform✓ Execution on XXX by code generation	<ul style="list-style-type: none">✓ Type-indexed types and functions✓ Staging + transform✓ Execution on XXX by code generation

- ▶ In the implementation we need “the best” of the two worlds
 - ▶ Type-indexed types from generic programming
 - ▶ Staged execution from deep embedding

Polytypic DSL

DSL	Monotypic (traditional)	Polytypic (data type generic)
Shallow embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Execution on the JVM	<ul style="list-style-type: none">✓ Type-indexed types and functions✓ Execution on the JVM
Deep embedding	<ul style="list-style-type: none">✓ Ordinary types and functions✓ Staging + transform✓ Execution on XXX by code generation	<ul style="list-style-type: none">✓ Type-indexed types and functions✓ Staging + transform✓ Execution on XXX by code generation <p>Polytypic Staging</p> <p>Nested Data Parallelism</p>

- ▶ In the implementation we need “the best” of the two worlds
 - ▶ Type-indexed types from generic programming
 - ▶ Staged execution from deep embedding

Framework: Polymorphic Embedding of DSLs

```
type Rep[A]    // abstract type constructor of representations
type PA[A] = Rep[PArray[A]]
trait PArray[A] {    // parallel array (to express parallelism)
  def length: Rep[Int]
  def map[R](f: Rep[A] => Rep[R]): PA[R]
  def zip[B](b: PA[B]): PA[(A,B)]
  ...
}
```

```
type Vector = PArray[Float] // parallel array
def dotProduct(vec1: Rep[Vector], vec2: Rep[Vector]): Rep[Float] =
  sum((vec1 zip vec2) map { case Pair(v1,v2) => v1 * v2 })
```

Framework: Polymorphic Embedding of DSLs

```
type Rep[A] // abstract type constructor of representations
type PA[A] = Rep[PArray[A]]
trait PArray[A] { // parallel array (to express parallelism)
  def length: Rep[Int]
  def map[R](f: Rep[A] => Rep[R]): PA[R]
  def zip[B](b: PA[B]): PA[(A,B)]
  ...
}
```

- The same code
- Two implementations
- Equivalent semantics

```
type Vector = PArray[Float] // parallel array
def dotProduct(vec1: Rep[Vector], vec2: Rep[Vector]): Rep[Float] =
  sum((vec1 zip vec2) map { case Pair(v1,v2) => v1 * v2 })
```

Shallow Embedding

```
type Rep[A] = A
```

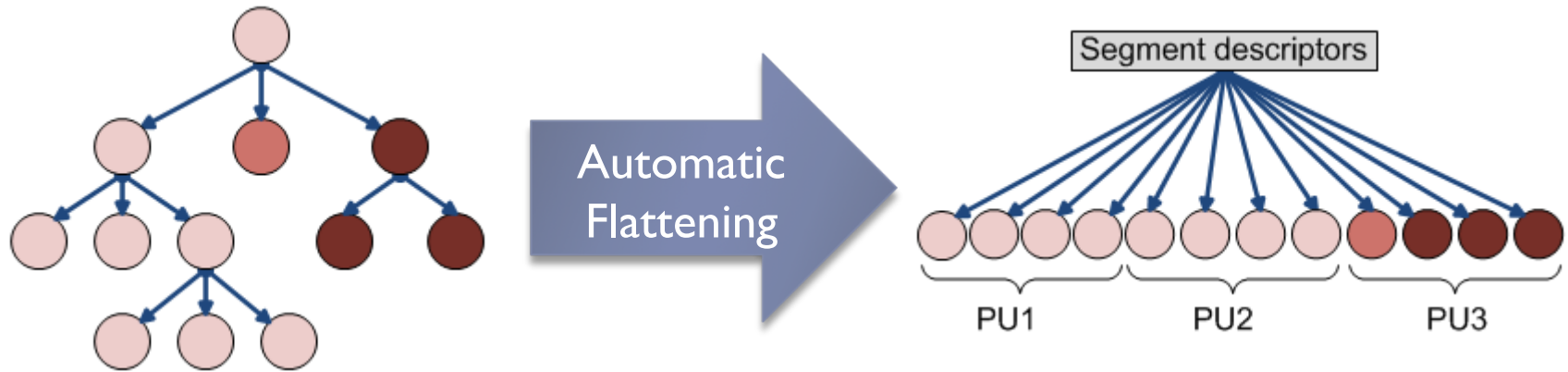
✓ Should be simple, good for testing and debugging

Deep Embedding

```
type Rep[A] = Exp[A]
```

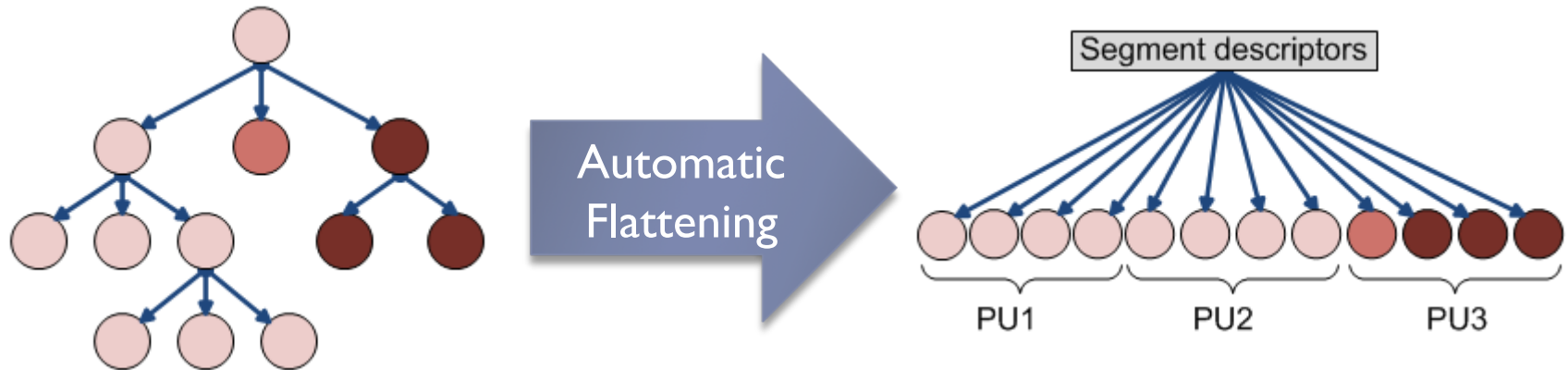
✓ Should generate an efficient code

The Key Idea – flattening transformation



- ▶ The data might be **irregular**
- ▶ **ill-balanced and not very parallel** at top level
- ▶ The one we want to write
- ▶ **regular** after flattening
- ▶ **Balanced** chunking
- ▶ The one we want to run

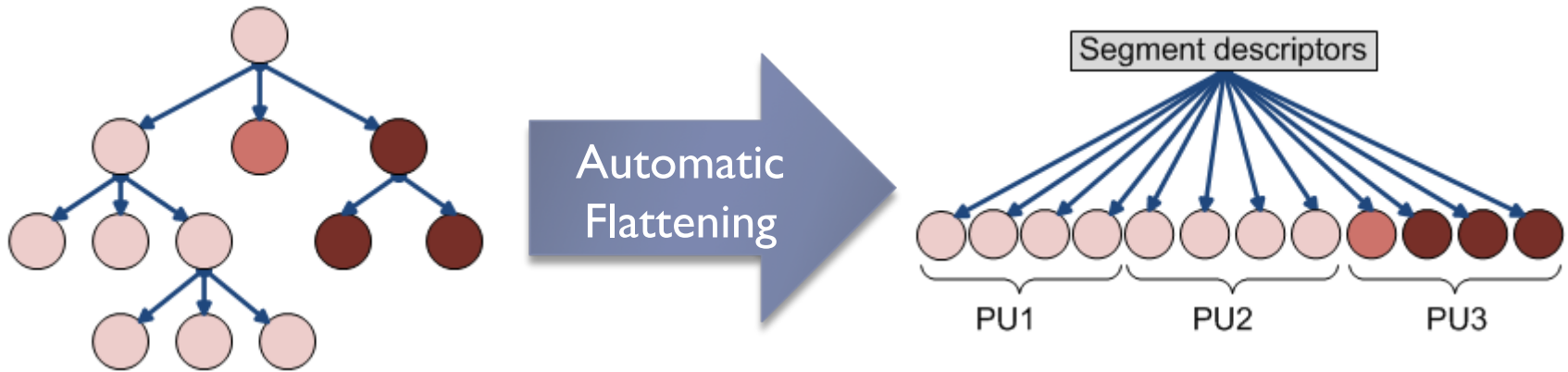
The Key Idea – flattening transformation



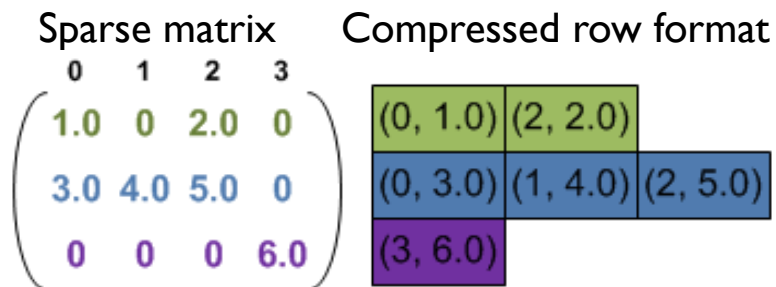
- ▶ The data might be **irregular**
- ▶ **ill-balanced and not very parallel** at top level
- ▶ The one we want to write
- ▶ **regular** after flattening
- ▶ **Balanced** chunking
- ▶ The one we want to run

Sparse matrix	Compressed row format									
$\begin{pmatrix} 0 & 1 & 2 & 3 \\ 1.0 & 0 & 2.0 & 0 \\ 3.0 & 4.0 & 5.0 & 0 \\ 0 & 0 & 0 & 6.0 \end{pmatrix}$	<table border="1"><tr><td>(0, 1.0)</td><td>(2, 2.0)</td><td></td></tr><tr><td>(0, 3.0)</td><td>(1, 4.0)</td><td>(2, 5.0)</td></tr><tr><td>(3, 6.0)</td><td></td><td></td></tr></table>	(0, 1.0)	(2, 2.0)		(0, 3.0)	(1, 4.0)	(2, 5.0)	(3, 6.0)		
(0, 1.0)	(2, 2.0)									
(0, 3.0)	(1, 4.0)	(2, 5.0)								
(3, 6.0)										

The Key Idea – flattening transformation

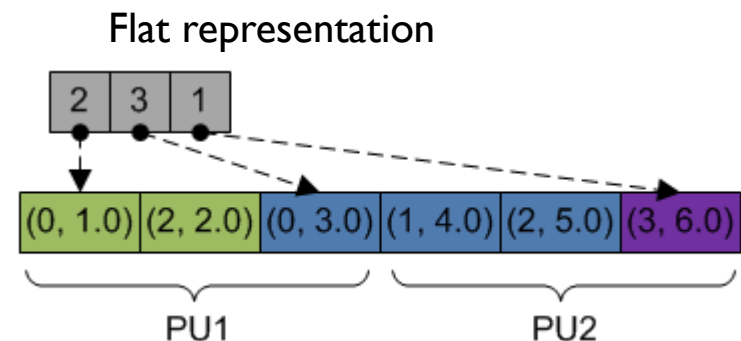


- ▶ The data might be **irregular**
- ▶ **ill-balanced and not very parallel** at top level
- ▶ The one we want to write



- ▶ **regular** after flattening
- ▶ **Balanced** chunking

- ▶ The one we want to run



How flattening happens

Nested Code	Flattened Code
<pre>p: A => B // primitive</pre>	<pre>type PA[A] = PArray[A] p^: PA[A] => PA[B] // p-lifted</pre>
<pre>def g(as: PA[A]) = as <u>map p</u> →</pre>	<pre>def g(as: PA[A]) = <u>p^</u>(as)</pre>
<pre>def h(m: PA[PA[A]]) = m map g →</pre>	<pre>def h(m: PA[PA[A]]) = m map g = m <u>map p^</u> // inline g = <u>p^^</u>(m) // ???</pre>

How flattening happens

Nested Code	Flattened Code
<pre>p: A => B // primitive</pre>	<pre>type PA[A] = PArray[A] p^: PA[A] => PA[B] // p-lifted</pre>
<pre>def g(as: PA[A]) = as <u>map</u> p</pre>	<pre>def g(as: PA[A]) = <u>p^</u>(as)</pre>
<pre>def h(m: PA[PA[A]]) = m map g</pre>	<pre>def h(m: PA[PA[A]]) = m map g = m <u>map</u> p^ // inline g = <u>p^^</u>(m) // ???</pre>

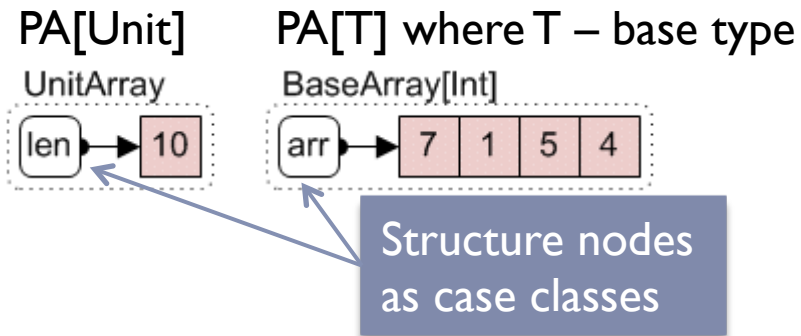
The key insight (we don't need p^^)

```
def p^^(m: PA[PA[A]]): PA[PA[A]] =  
  unconcat(m, p^(concat(m)))
```

```
def concat[A](nested: PA[PA[A]]): PA[A]
```

```
def unconcat[A,B](shape: PA[PA[A]], values: PA[B]): PA[PA[B]]
```

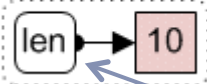
Data structures that support flattening



Data structures that support flattening

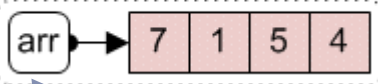
PA[Unit]

UnitArray



PA[T] where T – base type

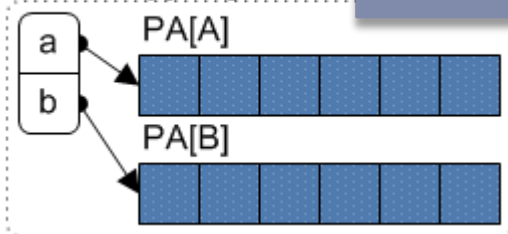
BaseArray[Int]



Structure nodes
as case classes

PA[(A,B)]

PairArray[A,B]



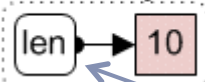
Constant time operations

```
def zip[A,B](a:PA[A], b:PA[B]):PA[(A,B)]  
  = PairArray(a, b)
```

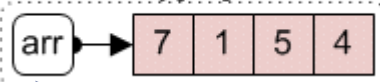
Data structures that support flattening

PA[Unit] PA[T] where T – base type

UnitArray



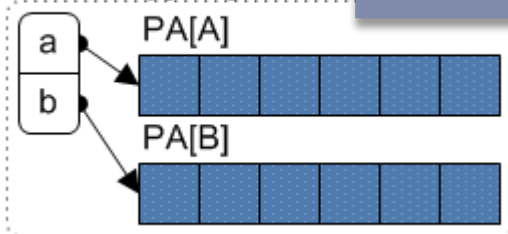
BaseArray[Int]



Structure nodes
as case classes

PA[(A,B)]

PairArray[A,B]

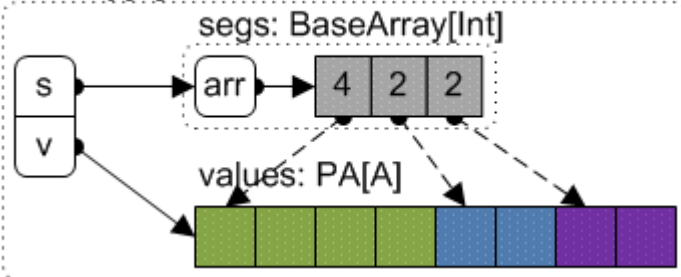


Constant time operations

```
def zip[A,B](a:PA[A], b:PA[B]):PA[(A,B)]  
  = PairArray(a, b)
```

PA[PA[A]]

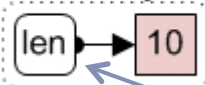
NArray[A]



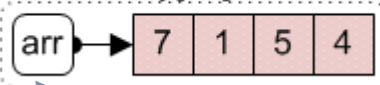
Data structures that support flattening

PA[Unit] PA[T] where T – base type

UnitArray



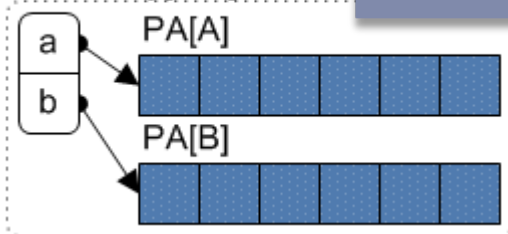
BaseArray[Int]



Structure nodes
as case classes

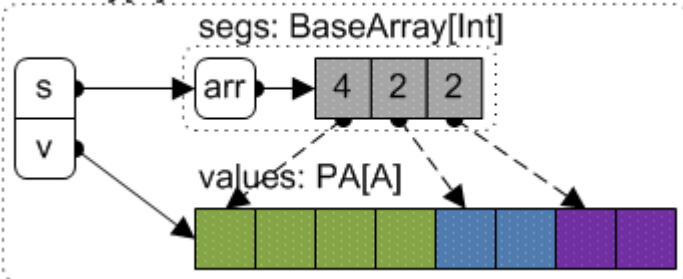
PA[(A,B)]

PairArray[A,B]



PA[PA[A]]

NArray[A]



Constant time operations

```
def zip[A,B](a:PA[A], b:PA[B]):PA[(A,B)]  
  = PairArray(a, b)
```

```
def concat[A](na: PA[PA[A]]): PA[A] =  
  na match { case NArray(vs, _) => vs }
```

```
def unconcat[A,B](shape:PA[PA[A]],  
                  vs:PA[B]): PA[PA[B]] =  
  shape match {  
    case NArray(_,segs) =>  
      NArray(vs,segs)  
  }
```

```
def p^^(m: PA[PA[A]]): PA[PA[A]] =  
  unconcat(m, p^(concat(m)))
```

Example (application specific types)

In the DSL we can construct types

```
T = Unit | Int | Float | Boolean // base types
  | (T1,T2) // product of types
  | (T1+T2) // sum of types
  | PArray[T] // nested array
```

$$\begin{matrix} & & M & & & & V & & = & R \\ \begin{pmatrix} (0, 1.0) & (2, 2.0) \\ (0, 3.0) & (1, 4.0) & (2, 5.0) \\ (3, 6.0) \end{pmatrix} & * & \begin{pmatrix} 1.0 & 0 & 2.0 & 0 \\ 3.0 & 4.0 & 5.0 & 0 \\ 0 & 0 & 0 & 6.0 \end{pmatrix} & * & \begin{pmatrix} 1.0 \\ 2.0 \\ 3.0 \\ 4.0 \end{pmatrix} & = & \begin{pmatrix} 7.0 \\ 26 \\ 24 \end{pmatrix} \end{matrix}$$

```
// Matrix in compressed row format
type SVector = PArray[(Int,Float)] // parallel array of products
type SMatrix = PArray[SVector] // nested array of rows
type Vector = PArray[Float] // dense vector
```

Example (Sparse Matrix Vector Multiplication)


$$\begin{matrix} & \text{M} & & \text{V} & = & \text{R} \\ \begin{pmatrix} (0, 1.0) & (2, 2.0) \\ (0, 3.0) & (1, 4.0) & (2, 5.0) \\ (3, 6.0) \end{pmatrix} & \begin{pmatrix} 1.0 & 0 & 2.0 & 0 \\ 3.0 & 4.0 & 5.0 & 0 \\ 0 & 0 & 0 & 6.0 \end{pmatrix} & * & \begin{pmatrix} 1.0 \\ 2.0 \\ 3.0 \\ 4.0 \end{pmatrix} & = & \begin{pmatrix} 7.0 \\ 26 \\ 24 \end{pmatrix} \end{matrix}$$

```
// Matrix in compressed row format
type SVector = PArray[(Int,Float)] // parallel array of products
type SMatrix = PArray[SVector]     // nested array of rows
type Vector = PArray[Float]       // dense vector
```

```
def sparseVectorMul(sv: SVector, vec: Vector): Float =
  sum(sv map { case (i,v) => vec(i) * v })
```

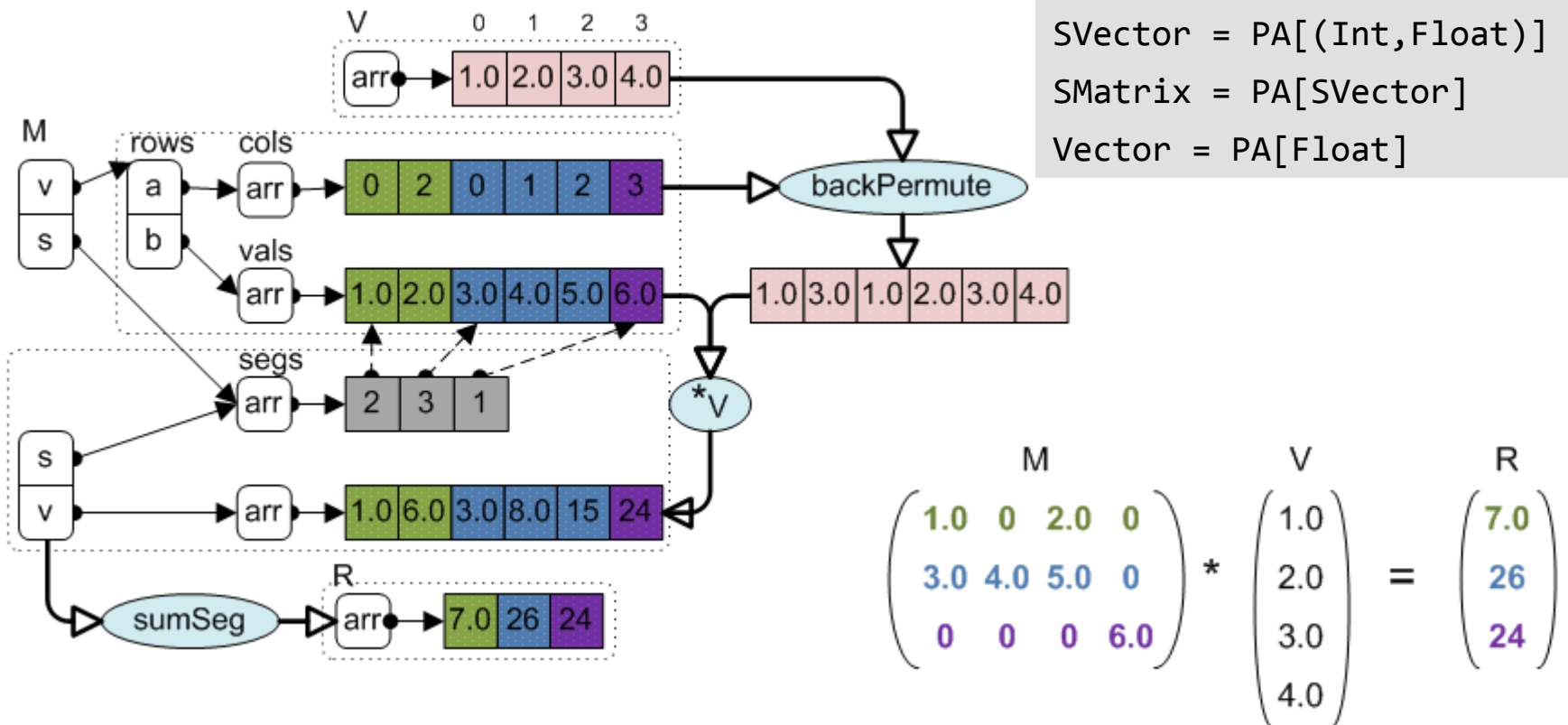
Inner parallelism

```
def smvm(matr: SMatrix, vec: Vector): Vector =
  for (row <- matr)
  yield sparseVectorMul(row, vec)
```

Outer parallelism

SMVM vectorized

```
def sparseVectorMul(sv: SVector, vec: Vector): Float =
  sum(sv map { case (i,v) => vec(i) * v })
def matrixVectorMul(matr: SMatrix, vec: Vector) =
  for (row <- matr) yield sparseVectorMul(row, vec)
```



Polytypic Staging

- ▶ Uses generic programming to capture domain semantics
- ▶ Allows the flattening of the DSL code by staged execution
- ▶ Is based on practical approaches: Scala-virtualized compiler, Polymorphic Embedding and LMS
- ▶ Not limited to NDP domain

Conclusions

- ▶ Nested data parallelism can be implemented in Scala as an embedded polytypic DSL
- ▶ To support flattening we need both staging and type-indexed data types
- ▶ Lightweight Polytypic Staging (LPS) is a framework for embedding of polytypic DSLs
- ▶ Nested Data Parallelism is a “killer app” for LPS
- ▶ WANTED: other polytypic domains?

References

- [1] Christian Hofer, Klaus Ostermann, Tillmann Rendel, and Adriaan Moors. *Polymorphic embedding of DSLs*. GPCE '08
- [2] Tiark Rompf, Martin Odersky. *Lightweight modular staging: a pragmatic approach to runtime code generation and compiled DSLs*. GPCE'10
- [3] Bruno C.d.S. Oliveira, Jeremy Gibbons. *Scala for generic programmers*. WGP'08

Q&A

???