

## Part XVII

# Staircase Join—Tree-Aware Relational (X)Query Processing

## Outline of this part

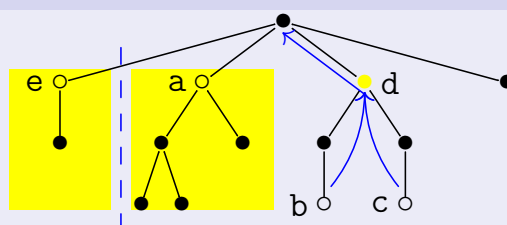
- 1 XPath Accelerator—Tree aware relational XML representation
  - Enhancing Tree Awareness
- 2 Staircase Join
  - Tree Awareness
  - Context Sequence Pruning
  - Staircases
- 3 Injecting ↗ into PostgreSQL
- 4 Outlook: More on Performance Tuning in MonetDB/XQuery

## Enhancing tree awareness

- We now know that the XPath Accelerator is a true **isomorphism** with respect to the XML skeleton **tree structure**.
  - Witnessed by our discussion of **shredder** ( $\mathcal{E}$ ) and **serializer** ( $\mathcal{E}^{-1}$ ).
- We will now see how the database kernel can benefit from a more elaborate **tree awareness** (beyond document order and semantics of the four major XPath axes).
- This will lead to the design of **staircase join**  $\sqsupset$ , the core of MonetDB/XQuery's XPath engine.
  - We will also discuss issues of how to tune  $\sqsupset$  to get the most out of modern CPUs and memory architectures.

## Tree awareness?

Document order and XPath semantics aside, what are further **tree properties** of value to a relational XML processor?



- 1 The **size of the subtree** rooted in node a is 5
- 2 The leaf-to-root **paths** of nodes b, c **meet** in node d
- 3 The **subtrees** rooted in e and a are necessarily **disjoint**

## Tree awareness ①: Subtree size

We have seen that tree property **subtree size** (① on previous slide) is implicitly present in a *pre/post*-based tree encoding:

$$post(v) - pre(v) = size(v) - level(v)$$

- To exploit property **subtree size**, we were able to find a means on the **SQL language level**, *i.e.*, **outside the database kernel**.
- ⇒ This led to *window shrink-wrapping* for the XPath descendant axis.

## Tree awareness on the SQL level

### Shrink-wrapping for the descendant axis

$$Q \equiv (c)/following::node()/descendant::node()$$

### *path(Q)*

```

SELECT  DISTINCT  $v_2.pre$ 
FROM    accel  $v_1$ , accel  $v_2$ 
WHERE    $v_1.pre > c.pre$ 
        AND  $v_1.pre < v_2.pre$ 
        AND  $v_1.post > c.post$ 
        AND  $v_1.post > v_2.post$ 
        AND  $v_2.pre \leq v_1.post + h$  AND  $v_2.post \geq v_1.pre + h$ 
ORDER BY  $v_2.pre$ 

```

## Tree awareness ②: Meeting ancestor paths

- Evaluation of axis `ancestor` can clearly benefit from knowledge about the exact element node where several given **node-to-root paths meet**.
  - For example:  
For context nodes  $c_1, \dots, c_n$ , determine their **lowest common ancestor**  $v = lca(c_1, \dots, c_n)$ .  
 $\Rightarrow$  Above  $v$ , produce result nodes once only.  
(This still produces duplicate nodes below  $v$ .)
- This knowledge *is* present in the encoding but is **not as easily expressed on the level of commonly available relational query languages** (such as, SQL or relational algebra).

## Flashback: XPath: Ensuring order is not for free

The strict XPath requirement to construct a result in document order may imply **sorting effort** depending on the actual XPath implementation strategy used by the processor.

```
(<x>
  <x><y id="0"/></x>
  <y id="1"/>
</x>)/descendant-or-self::x/child::y
```

$\Rightarrow$

```
(<y id="0"/>,
<y id="1"/>)
```

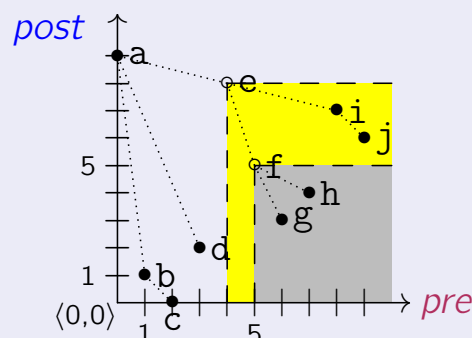
- In many implementations, the `descendant-or-self::x` step will yield the context node sequence  $(\langle x \rangle \dots \langle /x \rangle, \langle x \rangle \dots \langle /x \rangle)$  for the `child::y` step.
- Such implementations thus will typically extract `<y id="1"/>` before `<y id="0"/>` from the input document.

Flashback:  $(e, f)/\text{descendant}::\text{node}()$ 

## Context &amp; frag. encodings

context		
pre	post	...
5	5	
4	8	

accel		
pre	post	...
0	9	
1	1	
2	0	
3	2	
4	8	
5	5	
6	3	
7	4	
8	7	
9	6	

SQL query with expanded  $\text{window}()$  predicate

```

SELECT    DISTINCT v1.*
FROM      context v, accel v1
WHERE     v1.pre > v.pre AND v1.post < v.post
ORDER BY v1.pre

```

## Tree awareness ③: Disjoint subtrees

- An XPath location step  $cs/\alpha$  is evaluated for a context node **sequence**  $cs$ .
  - This “*set-at-a-time*” processing mode is key to the efficient evaluation of queries against bulk data. We want to map this into **set-oriented operations** on the RDBMS. (Remember: location step is translated into **join** between context node sequence and document encoding table *accel*.)
- But: If two context nodes  $c_{i,j} \in cs$  are in  $\alpha$ -relationship, **duplicates** and **out-of-order** results may occur.
  - Need efficient way to identify the  $c_i \in cs$  which are *not* in  $\alpha$ -relationship with any other  $c_j$  (for  $\alpha = \text{descendant}$ : “ $c_{i,j}$  in disjoint subtrees?”).

## Staircase Join: An injection of tree awareness

Since we fail to explain tree properties ② and ③ at the relational language level interface, we opt to **invade the database kernel** in a controlled fashion.<sup>46</sup>

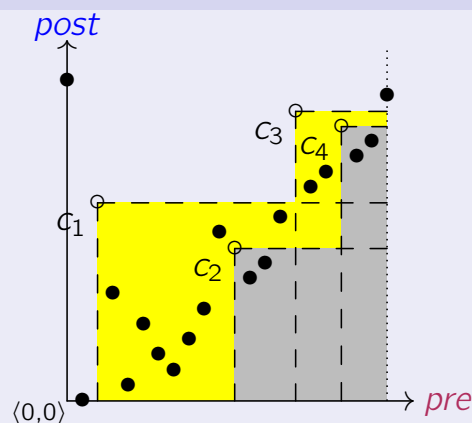
- Inject a new relational operator, **staircase join**  $\sqcup$ , into the relational query engine.
- Query translation and optimization in the presence of  $\sqcup$  continues to work like before (e.g., selection pushdown).
- The  $\sqcup$  algorithm encapsulates the necessary tree knowledge.  $\sqcup$  is a **local change** to the database kernel.

<sup>46</sup>Remember: All of this is optional. XPath Accelerator is a purely relational XML document encoding, working on top of *any* RDBMS.

## Tree awareness: Window overlap, coverage

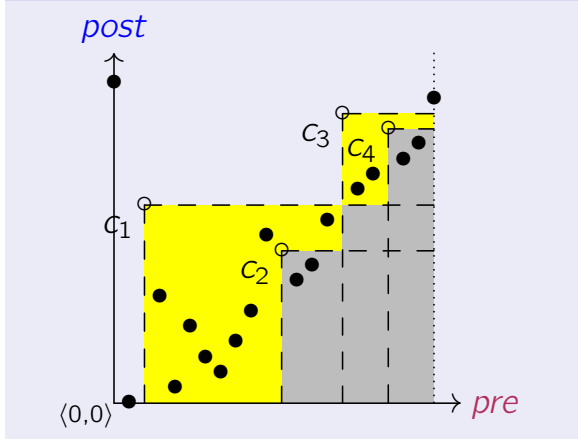
Location step  $(c_1, c_2, c_3, c_4)/\text{descendant}::\text{node}()$ . The pairs  $(c_1, c_2)$  and  $(c_3, c_4)$  are in descendant-relationship:

### Window overlap and coverage (descendant axis)

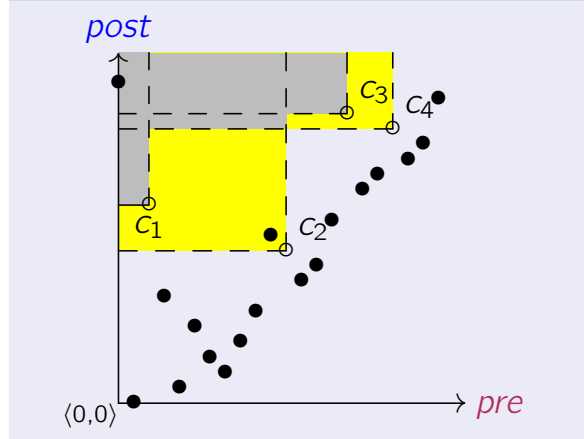


# Tree awareness: Window overlap, coverage

Axis window overlap (descendant axis)

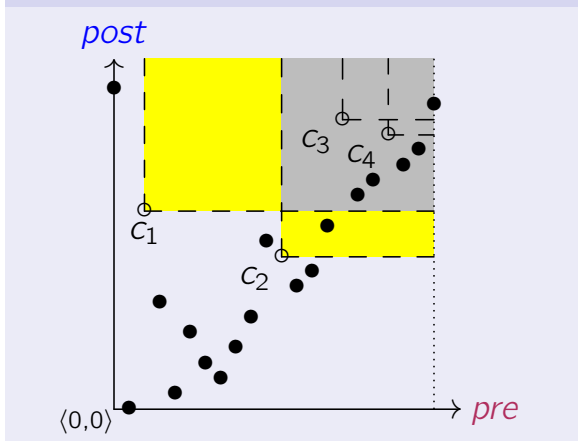


Axis window overlap (ancestor axis)

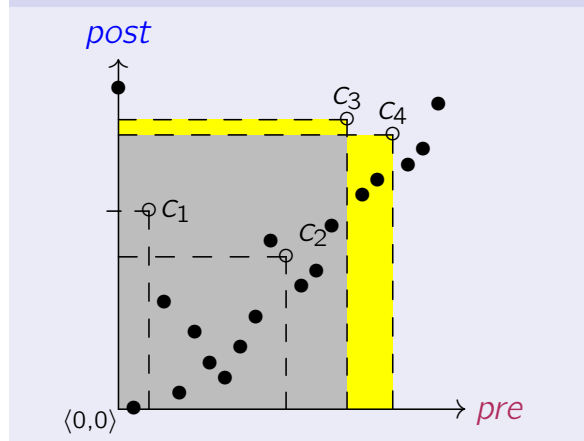


# Tree awareness: Window overlap, coverage

Axis window overlap (following axis)



Axis window overlap (preceding axis)



## Context node sequence pruning

We can turn these observations about axis window overlap and coverage into a simple strategy to **prune the initial context node sequence** for an XPath location step.

### Context node sequence pruning

Given  $cs/\alpha$ , determine minimal  $cs^- \subseteq cs$ , such that

$$cs/\alpha = cs^-/\alpha .$$

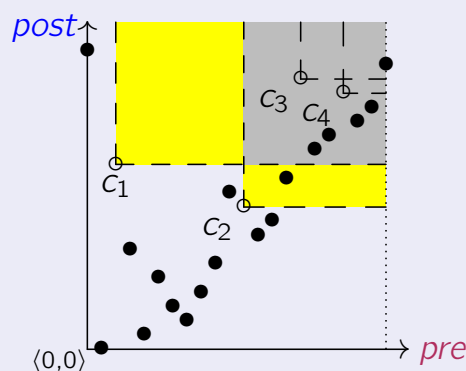
We will see that this minimization leads to axis step evaluation on the *pre/post* plane, which never emits duplicate nodes or out-of-order results.<sup>47</sup>

<sup>47</sup>The ancestor axis needs a bit more work here.

## Context node pruning: following axis

Once context pruning for the following axis is complete, all remaining context nodes relate to each other on the ancestor/descendant axes:

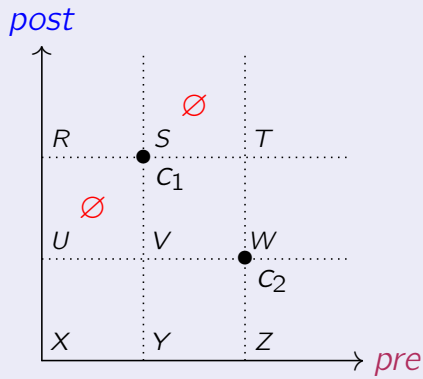
### Covering nodes $c_{1,2}$ in descendant relationship





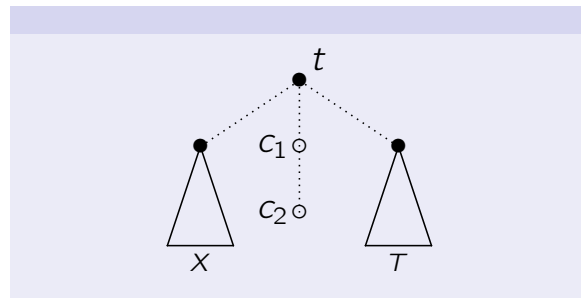
# Empty regions in the *pre/post* plane

## Relating two context nodes ( $c_1, c_2$ ) on the plane



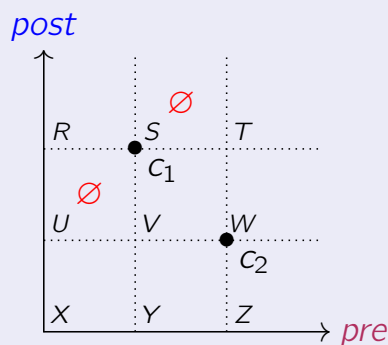
## Empty regions?

Given  $c_{1,2}$  on the left, why are the regions  $U, S$  marked  $\emptyset$  guaranteed *to not hold any nodes*?



# Context pruning (following axis)

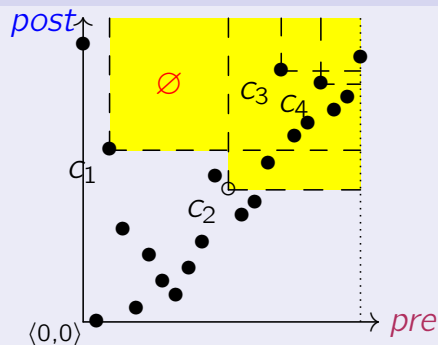
## $(c_1, c_2)/following::node()$



$$\begin{aligned}
 (c_1, c_2)/following::node() &\equiv SUTUW \\
 &\equiv T UW \\
 &\equiv (c_2)/following::node()
 \end{aligned}$$

## Context pruning (following axis)

### Context pruning (following axis)



### Context pruning (following axis)

Replace context node sequence  $cs$  by singleton sequence  $(c)$ ,  $c \in cs$ , with  $post(c)$  minimal.

## Context pruning (preceding axis)

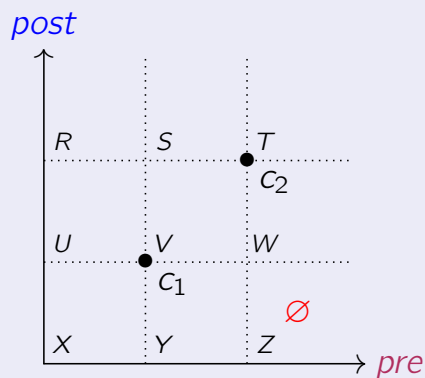
### Context pruning (preceding axis)

Replace context node sequence  $cs$  by singleton sequence  $(c)$ ,  $c \in cs$ , with  $pre(c)$  maximal.

- Regardless of initial context size, axes following and preceding yield simple **single region queries**.
- We focus on descendant and ancestor now.

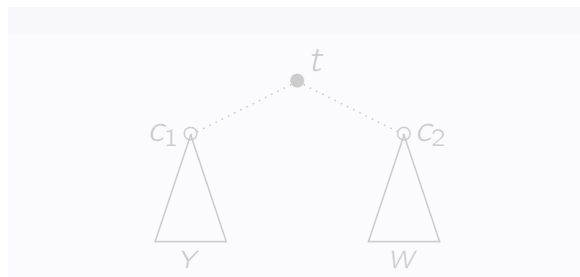
## More empty regions

Remaining context nodes  
 $c_1, c_2$  after pruning for  
descendant axis



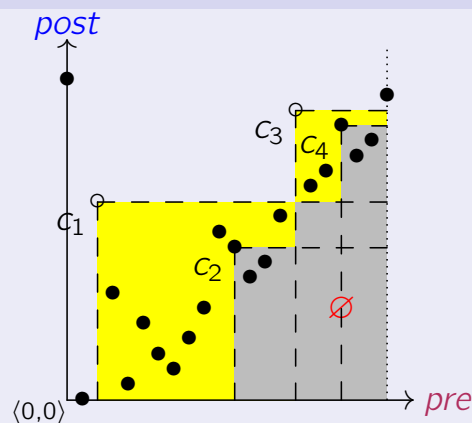
Empty region?

Why is region  $Z$  marked  $\emptyset$   
guaranteed to be empty?



## Context pruning (descendant axis)

Context pruning (descendant axis)



- The region marked  $\emptyset$  above is a region of type  $Z$  (previous slide). In general, a **non-singleton sequence remains**.

## Context pre-processing: Pruning

```
prune_contextdesc(context : TABLE(pre, post))
```

```

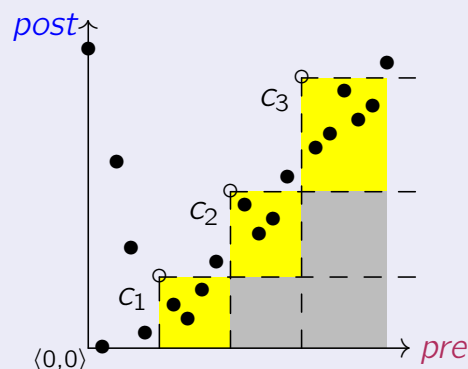
begin
  result ← CREATE TABLE(pre, post);
  prev ← 0;
  foreach c in context do
    /* retain node only if post rank increases */
    if c.post > prev then
      APPEND c TO result;
      prev ← c.post;
    /* return new context table */
  return result;
end

```

## “Staircases” in the *pre/post* plane

Note that after context pruning, the remaining context nodes form a **proper “staircase”** in the plane. (This is an important assumption in the following.)

### Context pruning & “staircase”



## Flashback: Intersecting ancestor paths

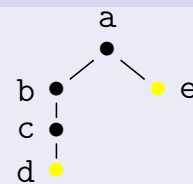
Even with pruning applied, duplicates and out-of-order results may still be generated due to **intersecting ancestor paths**.

- We have observed this before: apply function `ancestors( $c_1, c_2$ )` where  $c_1$  ( $c_2$ ) denotes the element node with tag `d` (`e`) in the sample tree below.  
(Nodes  $c_{1,2}$  would *not* have been removed during pruning.)

### Simulate XPath ancestor via parent axis

```
declare function
  ancestors($n as node()) as node()*
{ if (fn:empty($n)) then ()
  else (ancestors($n/..), $n/..)
}
```

### Sample tree



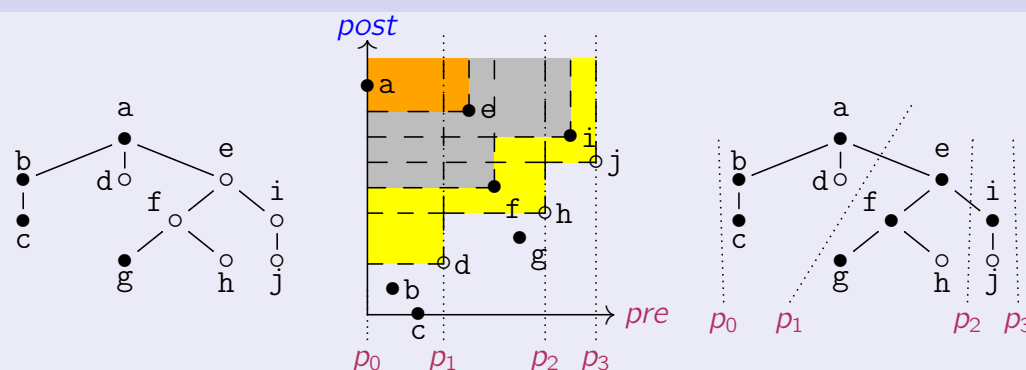
**Remember:** `ancestors((d,e))` yielded `(a,b,a,c)`.

## Separation of ancestor paths

**Idea:** try to **separate** the ancestor paths by defining suitable **cuts** in the XML fragment tree.

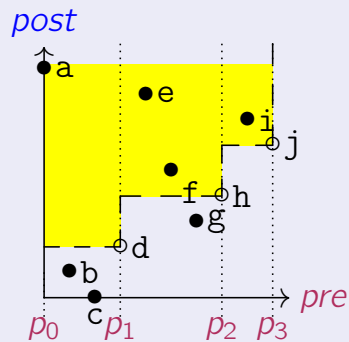
- Stop node-to-root traversal if a cut is encountered.

### Path separation (ancestor axis)



## Parallel scan along the *pre* dimension

### Separating ancestor paths



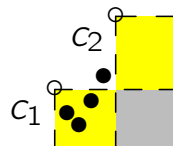
Scan partitions (intervals):  $[p_0, p_1)$ ,  $[p_1, p_2)$ ,  $[p_2, p_3)$ .

- Can scan in **parallel**. Partition results may be concatenated.
- Context pruning reduces numbers of partitions to scan.

## Basic Staircase Join (descendant)

```
⌊desc(accel : TABLE(pre, post), context : TABLE(pre, post))
```

```
begin
  result ← CREATE TABLE(pre, post);
  foreach successive pair (c1, c2) in context do
    scanpartition(c1.pre + 1, c2.pre - 1, c1.post, <);
  c ← last node in context;
  n ← last node in accel;
  scanpartition(c.pre + 1, n.pre, c.post, <);
  return result;
end
```



## Partition scan (sub-routine)

*scanpartition*(*pre*<sub>1</sub>, *pre*<sub>2</sub>, *post*,  $\theta$ )

```

begin
  | for i from pre1 to pre2 do
  |   | if accel[i].post  $\theta$  post then
  |   |   | APPEND accel[i] TO result;
  |   |
  |   end
end

```

Notation *accel*[*i*] does not imply random access to document encoding:

- Access is strictly **forward sequential** (also *between* invocations of *scanpartition*(.)).

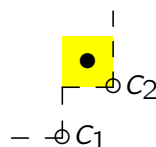
## Basic Staircase Join (ancestor)

$\sqcup_{\text{anc}}$ (*accel* : TABLE(*pre*, *post*), *context* : TABLE(*pre*, *post*))

```

begin
  | result  $\leftarrow$  CREATE TABLE(pre, post);
  | c  $\leftarrow$  first node in context;
  | n  $\leftarrow$  first node in accel;
  | scanpartition(n.pre, c.pre - 1, c.post, >);
  | foreach successive pair (c1, c2) in context do
  |   | scanpartition(c1.pre + 1, c2.pre - 1, c2.post, >);
  |   return result;
end

```



## Basic Staircase Join: Summary

- The operation of **staircase join** is perhaps most closely described as **merge join** with a **dynamic range predicate**: the join predicate traces the staircase boundary:
  - $\sqcup$  scans the *accel* and *context* tables and populates the *result* table sequentially in document order,
  - $\sqcup$  scans both tables once for an entire context sequence,
  - $\sqcup$  never delivers duplicate nodes.
- $\sqcup$  works correctly only if *prune\_context*(·) has previously been applied.
  - *prune\_context*(·) may be **inlined** into  $\sqcup$ , thus performing context pruning *on-the-fly*.



## Pruning on-the-fly

```
 $\sqcup_{desc}(accel:TABLE(pre, post), context:TABLE(pre, post))$ 
```

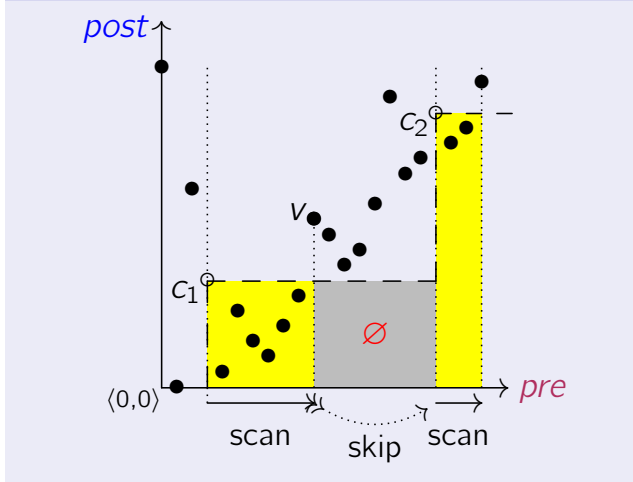
```

begin
  result ← CREATE TABLE(pre, post);
  c1 ← first node in context;
  while (c2 ← next node in context) do
    if c2.post < c1.post then
      /* prune */
    else
      scanpartition(c1.pre + 1, c2.pre - 1, c1.post, <);
      c1 ← c2;
  c ← last node in context;
  n ← last node in accel;
  scanpartition(c.pre + 1, n.pre, c.post, <);
  return result;
end

```



## Skip ahead, if possible

 $(c_1, c_2)/\text{descendant}::\text{node}()$ 

- While scanning the partition associated with  $c_{1,2}$ :
  - $v$  is outside staircase boundary, thus not part of the result.
  - No node beyond  $v$  in result ( $\emptyset$ -region of type  $Z$ ).
- ⇒ Can **terminate scan early and skip ahead** to  $pre(c_2)$ .

## Skipping for the descendant axis

 $\text{scanpartition}_{\text{desc}}(\text{pre}_1, \text{pre}_2, \text{post})$ 

```

begin
  for  $i$  from  $\text{pre}_1$  to  $\text{pre}_2$  do
    if  $\text{accel}[i].\text{post} < \text{post}$  then
      APPEND  $\text{accel}[i]$  TO result;
    else
      /* on the first offside node, terminate scan */ break;
  end
end

```

**Note:** keyword **break** transfers control out of innermost enclosing loop (cf. C, Java).

## Effectiveness of skipping

- Enable skipping in  $scanpartition(\cdot)$ . Then, for each node in  $context$ , we either
  - ① hit a node to be copied into table  $result$ , or
  - ② encounter an offside node (node  $v$  on slide 1128) which leads to a skip to a known  $pre$  value ( $\rightarrow$  positional access).
- To produce the final result,  $\surd$  thus never touches more than

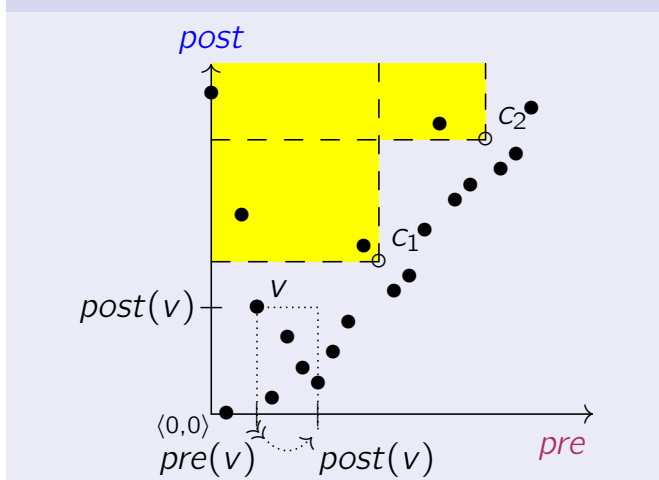
$$|context| + |result|$$

nodes in the plane (without skipping:  $|context| + |accel|$ ).

- In practice:  $> 90\%$  of nodes in table  $accel$  are skipped.

## Skipping for the ancestor axis

### Skipping over the subtree of $v$



Encounter  $v$  outside staircase boundary

$\Downarrow$

$v$  and subtree below  $v$  in preceding axis of context node.

### How far to skip?

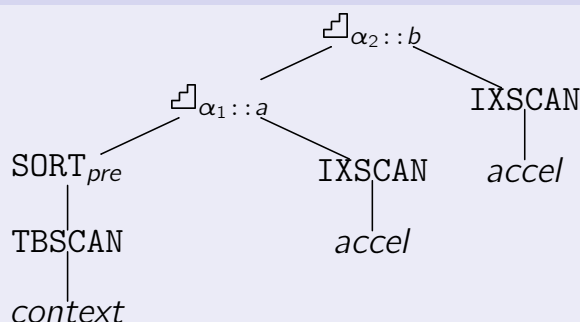
Conservative estimate:  
 $size(v) \geq post(v) - pre(v)$

## Injecting $\sqcup$ into PostgreSQL

**PostgreSQL** (<http://postgresql.org/>): Conventional disk-based RDBMS, SQL interface.

- Detection of  $\sqcup$  applicability on SQL level (self-join with conjunctive range selection on columns of type tree<sup>48</sup>).

### Algebraic query plan for two-step XPath location path



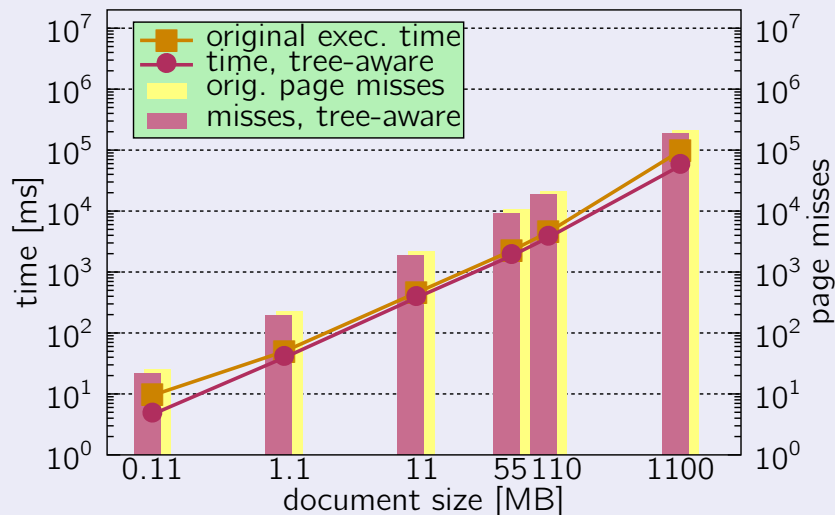
<sup>48</sup>PostgreSQL is highly extensible, also permits introduction of new column types.

## Injecting $\sqcup$ into PostgreSQL

- Create clustered ascending B-tree index on column *pre* of table *accel*.
  - Standard no-frills PostgreSQL B-tree index, entered with search predicates of the form  $pre \geq c.pre$  ( $c$  context node).
  - B-tree on column *pre* also used for **skipping**.
- Following performance figures obtained on a 2.2 GHz Dual Intel<sup>TM</sup> Pentium 4, 2 GB RAM, PostgreSQL 7.3.3.
  - Compares  $\sqcup$ -enabled (tree-aware) PostgreSQL with vanilla PostgreSQL instance.
  - Evaluate XPath location path  $/descendant::a/\alpha::b$  on document instances of up to 1.1 GB serialized size.

Injecting  $\alpha$  into PostgreSQL

/descendant::a/descendant::b

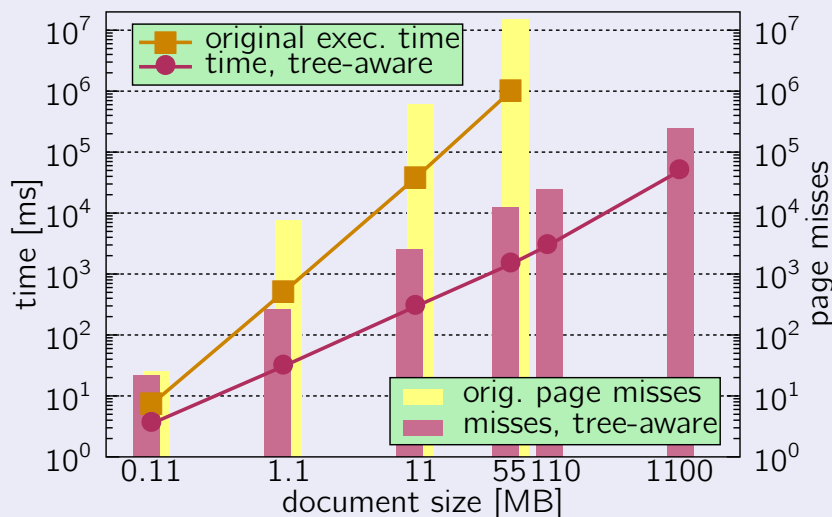
Injecting  $\alpha$  into PostgreSQL

**For  $\alpha = \text{descendant}$  observe:**

- For *both* PostgreSQL instances, query evaluation time grows **linearly** with the input XML document size (since the results size grows linearly).
- For the original instance, this is due to **window shrink-wrapping** (expressible at the SQL level).

Injecting  $\alpha$  into PostgreSQL

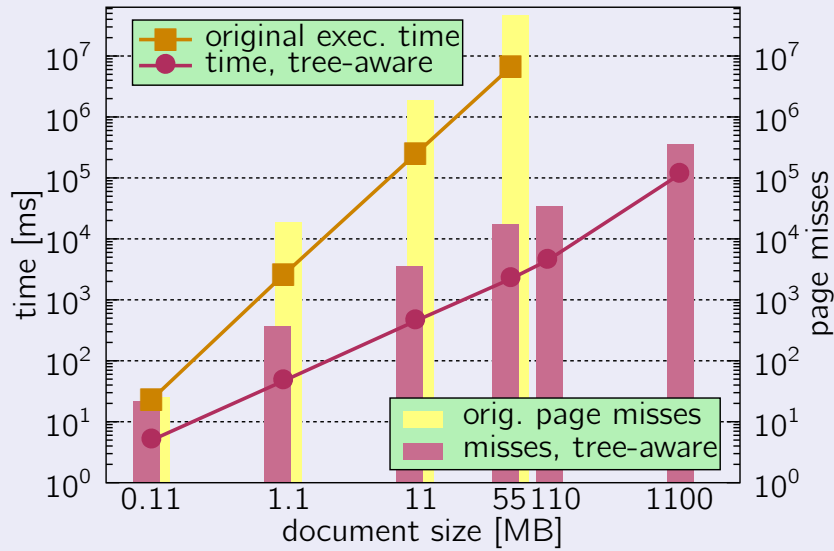
/descendant::a/ancestor::b

Injecting  $\alpha$  into PostgreSQL

- **For  $\alpha \in \{\text{ancestor, preceding, following}\}$  observe:**
  - For the  $\alpha$ -**enabled** PostgreSQL instance, query evaluation time grows **linearly** with the input XML document (and result) size.  
For the **original** instance, query evaluation time grows **quadratically** ( $| accel |$  scans of table *accel* performed).
  - Original instance is incapable of completing experiment in reasonable time ( $> 15$  mins for XML input size of 55 MB).
- **Generally:**
  - The number of **buffer page misses** (= necessary I/O operations) determines evaluation time.

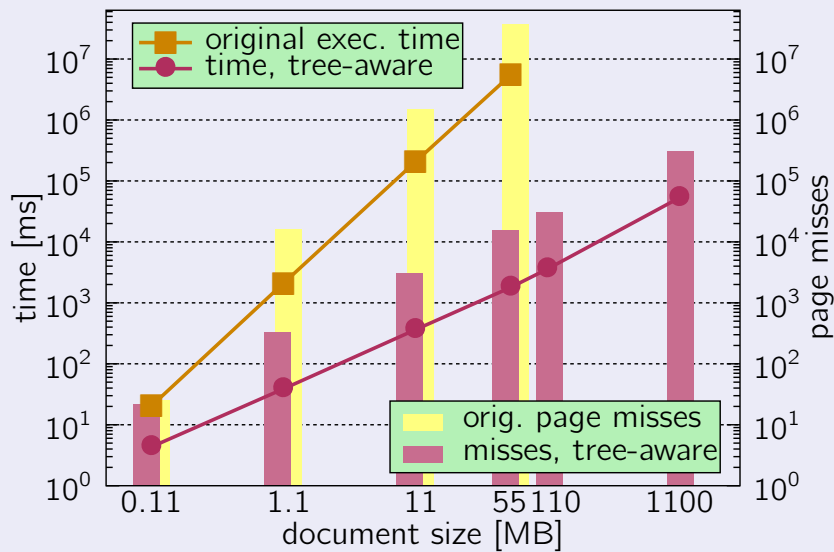
# Injecting $\sqcup$ into PostgreSQL

`/descendant::a/preceding::b`



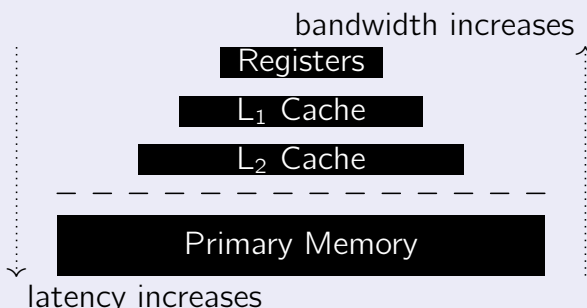
# Injecting $\sqcup$ into PostgreSQL

`/descendant::a/following::b`



# MonetDB/XQuery: Targetting modern CPU/memory architectures

## Memory Hierarchy



- Computation performed with CPU registers only.
- **Cache miss may escalate:**  $L_1 \rightarrow L_2 \rightarrow \text{RAM}$ , data transport all the way back:  $L_1 \leftarrow L_2 \leftarrow \text{RAM}$ .
- Data transport in **cache line granularity**.

# CPU/cache characteristics

## Intel™ Dual Pentium 4 (Xeon)<sup>49</sup>

### CPU/Cache Characteristics

Clock frequency		2.2 GHz
L <sub>1</sub> /L <sub>2</sub> cache size		8 kB/512 kB
L <sub>1</sub> /L <sub>2</sub> cache line size	$LS_{L_1}/LS_{L_2}$	32 byte/128 byte
L <sub>1</sub> miss latency	$L_{L_1}$	28 cycles $\hat{=}$ 12.7 ns
L <sub>2</sub> miss latency	$L_{L_2}$	387 cycles $\hat{=}$ 176 ns

- For this CPU, a **full cache miss** implies a **stall of the CPU** for  $28 + 387 = 415$  cycles (cy).



<sup>49</sup>Measure these characteristics for your CPU with Stefan Manegold's *Calibrator*, <http://monetdb.cwi.nl/Calibrator/>.

## Staircase join: Wrap-up

- Standard B<sup>+</sup>-tree implementation suffices to support  $\sqsupset$ .
  - A **single** B<sup>+</sup>-tree indexes the *pre/post* plane as well as the context node sequence.
    - ⇒ Less index pages compete for valuable buffer space.
- $\sqsupset$  derives pruning and skipping information from the plane itself, using **simple integer arithmetic and comparisons**.
  - Simple  $\sqsupset$  logic leads to **simple memory access pattern and control flow**.
    - ⇒ Branches in inner  $\sqsupset$  loops are highly predictable, facilitating **speculative execution** in the CPU.

Predictable branches?

Explain why!