

## Part V

# SAX—Simple API for XML

# Outline of this part

- 1 SAX Events
- 2 SAX Callbacks
- 3 SAX and the XML Tree Structure
- 4 SAX and Path Queries
  - Path Query Evaluation
- 5 Final Remarks on SAX

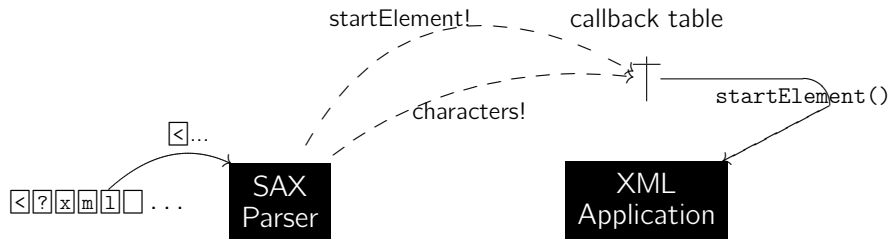
## SAX—Simple API for XML

- **SAX<sup>5</sup> (Simple API for XML)** is, unlike DOM, *not* a W3C standard, but has been developed jointly by members of the XML-DEV mailing list (*ca.* 1998).
- SAX processors use **constant space**, regardless of the XML input document size.
  - Communication between the SAX processor and the back-end XML application does *not* involve an intermediate tree data structure.
  - Instead, the **SAX parser sends events** to the application whenever a certain piece of XML text has been recognized (*i.e.*, parsed).
  - The **back-end acts on/ignores events** by populating a **callback function table**.

---

<sup>5</sup><http://www.saxproject.org/>

## Sketch of SAX's mode of operations



- A SAX processor reads its input document **sequentially** and **once** only.
- No memory of what the parser has seen so far is retained while parsing. As soon as a significant bit of XML text has been recognized, an **event** is sent.
- The application is able to act on events **in parallel** with the parsing progress.



# SAX Events

- To meet the constant memory space requirement, SAX reports **fine-grained parsing events** for a document:

Event	... reported when seen	Parameters sent
<i>startDocument</i>	<code>&lt;?xml...?&gt;</code> <sup>6</sup>	
<i>endDocument</i>	<code>&lt;EOF&gt;</code>	
<i>startElement</i>	<code>&lt;t a<sub>1</sub>=v<sub>1</sub> ... a<sub>n</sub>=v<sub>n</sub>&gt;</code>	<i>t</i> , ( <i>a</i> <sub>1</sub> , <i>v</i> <sub>1</sub> ), ..., ( <i>a</i> <sub><i>n</i></sub> , <i>v</i> <sub><i>n</i></sub> )
<i>endElement</i>	<code>&lt;/t&gt;</code>	<i>t</i>
<i>characters</i>	<i>text content</i>	Unicode buffer ptr, length
<i>comment</i>	<code>&lt;!--c--&gt;</code>	<i>c</i>
<i>processingInstruction</i>	<code>&lt;?t pi?&gt;</code>	<i>t</i> , <i>pi</i>
	⋮	

<sup>6</sup>**N.B.:** Event *startDocument* is sent even if the optional XML text declaration should be missing.

## dilbert.xml

```

1 <?xml encoding="utf-8"?> *1
2 <bubbles> *2
3   <!-- Dilbert looks stunned --> *3
4   <bubble speaker="phb" to="dilbert"> *4
5     Tell the truth, but do it in your usual engineering way
6     so that no one understands you. *5
7   </bubble> *6
8 </bubbles> *7 *8

```

Event <sup>7 8</sup>	Parameters sent
*1	<i>startDocument</i>
*2	<i>startElement</i> <i>t</i> = "bubbles"
*3	<i>comment</i> <i>c</i> = "_Dilbert looks stunned_"
*4	<i>startElement</i> <i>t</i> = "bubble", ("speaker","phb"), ("to","dilbert")
*5	<i>characters</i> <i>buf</i> = "Tell the...understands you.", <i>len</i> = 99
*6	<i>endElement</i> <i>t</i> = "bubble"
*7	<i>endElement</i> <i>t</i> = "bubbles"
*8	<i>endDocument</i>

<sup>7</sup>Events are reported in **document reading order** \*1, \*2, ..., \*8.

<sup>8</sup>**N.B.:** Some events suppressed (white space).

# SAX Callbacks

- To provide an efficient and tight **coupling** between the SAX **frontend** and the application **backend**, the SAX API employs **function callbacks**:<sup>9</sup>
  - Before parsing starts, the application **registers function references** in a table in which each event has its own slot:

Event	Callback		Event	Callback
⋮			⋮	
<i>startElement</i>	?	— →	<i>startElement</i>	<i>startElement ()</i>
<i>endElement</i>	?	<i>SAXregister(startElement,</i> <i>startElement ())</i>	<i>endElement</i>	<i>endElement ()</i>
⋮		<i>SAXregister(endElement,</i> <i>endElement ())</i>	⋮	
⋮			⋮	

- The application alone decides on the implementation of the functions it registers with the SAX parser.
- Reporting an event**  $\star_j$  then amounts to call the function (with parameters) registered in the appropriate table slot.

<sup>9</sup>Much like in event-based GUI libraries.



## Java SAX API

In Java, populating the callback table is done via implementation of the SAX `ContentHandler` interface: a `ContentHandler` object represents the callback table, its methods (e.g., `public void endDocument ()`) represent the table slots.

**Example:** Reimplement *content.cc* shown earlier for DOM (find all XML text nodes and print their content) using SAX (pseudo code):

*content (File f)*

```
// register the callback,
// we ignore all other events
SAXregister (characters, printText);
SAXparse (f);
return;
```

*printText ((Unicode) buf, Int len)*

```
Int i;
foreach  $i \in 1 \dots len$  do
   $\lfloor$  print (buf[i]);
return;
```



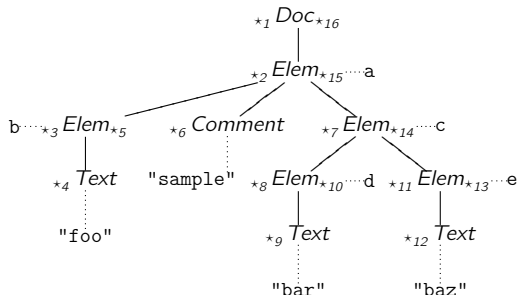
# SAX and the XML Tree Structure

- Looking closer, the **order** of SAX events reported for a document is determined by a **preorder traversal** of its document tree<sup>10</sup>:

Sample XML document

```

1  *1
2  <a>*2
3    <b>*3 foo*4 </b>*5
4    <!--sample-->*6
5    <c>*7
6      <d>*8 bar*9 </d>*10
7      <e>*11 baz*12 </e>*13
8    </c>*14
9  </a>*15 *16
  
```



**N.B.:** An *Elem* [*Doc*] node is associated with two SAX events, namely *startElement* and *endElement* [*startDocument*, *endDocument*].

<sup>10</sup>Sequences of sibling *Char* nodes have been collapsed into a single *Text* node.

## Challenge

- This **left-first depth-first** order of SAX events is well-defined, but appears to make it hard to answer certain queries about an XML document tree.

 Collect all direct children nodes of an *Elem* node.

In the example on the previous slide, suppose your application has just received the *startElement*(*t* = "a") event  $\star_2$  (i.e., the parser has just parsed the opening element tag <a>).

With the remaining events  $\star_3 \dots \star_{16}$  still to arrive, can your code detect all the immediate children of *Elem* node a (i.e., *Elem* nodes b and c as well as the *Comment* node)?

The previous question can be answered more generally:

*SAX events are sufficient to **rebuild the complete XML document tree** inside the application. (Even if we most likely don't want to.)*

## SAX-based tree rebuilding strategy (sketch):

- ① [*startDocument*]  
Initialize a **stack**  $S$  of **node IDs** (e.g.  $\in \mathbb{Z}$ ). **Push** first ID for this node.
- ② [*startElement*]  
Assign a **new ID** for this node. **Push** the ID onto  $S$ .<sup>11</sup>
- ③ [*characters, comment, ...*]  
Simply assign a new node ID.
- ④ [*endElement, endDocument*]  
**Pop**  $S$  (no new node created).

**Invariant:** The **top of**  $S$  holds the identifier of the current **parent node**.

<sup>11</sup>In callbacks ② and ③ we might wish to store further node details in a table or similar summary data structure.

## SAX Callbacks

**SAX callbacks** to rebuild XML document tree:

- We maintain a summary table of the form

ID	NodeType	Tag	Content	ParentID
----	----------	-----	---------	----------

- *insert* (*id*, *type*, *t*, *c*, *pid*) inserts a row into this table.
- Maintain stack *S* of node IDs, with operations *push(id)*, *pop()*, *top()*, and *empty()*.

*startDocument* ()

```

id ← 0;
S.empty();
insert (id, Doc, □, □, □);
S.push(id);
return;

```

*endDocument* ()

```

S.pop();
return;

```

# SAX Callbacks

startElement ( $t, (a_1, v_1), \dots$ )

```

id ← id + 1;
insert (id, Elem, t, □, S.top());
S.push(id);
return;

```

characters ( $buf, len$ )

```

id ← id + 1;
insert (id, Text, □, buf[1...len], S.top());
return;

```

endElement ( $t$ )

```

S.pop();
return;

```

comment ( $c$ )

```

id ← id + 1;
insert (id, Comment, □, c, S.top());
return;

```

- Run against the example given above, we end up with the following summary table:

ID	NodeType	Tag	Content	ParentID
0	<i>Doc</i>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
1	<i>Elem</i>	a	<input type="checkbox"/>	0
2	<i>Elem</i>	b	<input type="checkbox"/>	1
3	<i>Text</i>	<input type="checkbox"/>	"foo"	2
4	<i>Comment</i>	<input type="checkbox"/>	"sample"	1
5	<i>Elem</i>	c	<input type="checkbox"/>	1
6	<i>Elem</i>	d	<input type="checkbox"/>	5
7	<i>Text</i>	<input type="checkbox"/>	"bar"	6
8	<i>Elem</i>	e	<input type="checkbox"/>	5
9	<i>Text</i>	<input type="checkbox"/>	"baz"	8

- Since XML defines tree structures only, the ParentID column is all we need to recover the complete node hierarchy of the input document.

### Walking the XML node hierarchy?

Explain how we may use the summary table to find the (a) *children*, (b) *siblings*, (c) *ancestors*, (d) *descendants* of a given node (identified by its ID).

## Final Remarks on SAX

- For an XML document fragment shown on the left, SAX might actually report the events indicated on the right:

XML fragment		XML + SAX events	
1	<affiliation>	1	<affiliation>* <sub>1</sub>
2	AT&amp;T Labs	2	AT* <sub>2</sub> &amp;* <sub>3</sub> T Labs
3	</affiliation>	3	* <sub>4</sub> </affiliation>* <sub>5</sub>

---

\*<sub>1</sub> *startElement*(affiliation)  
 \*<sub>2</sub> *characters*(" \n \t AT", 5)  
 \*<sub>3</sub> *characters*("&", 1)  
 \*<sub>4</sub> *characters*("T \t Labs \n", 7)  
 \*<sub>5</sub> *endElement*(affiliation)



**White space** is reported.

**Multiple characters events** may be sent for text content (although adjacent).

(Often SAX parsers break text on entities, but may even report each character on its own.)