

Source Code Search

and other applications of in-memory technology in software development

Oleksandr Panchenko

29.06.2011

There is a huge amount of source code:

– at SAP:

- Business ByDesign and NetWeaver – 72 Mio ABAP statements, 124k development objects (classes, reports, function modules), 1,6 Mio includes
- NetWeaver Basis – 12k database tables, 800k SQL statements
- Business Suite – 400 Mio lines of code

– at Microsoft:

- Microsoft Vista – 120 Mio lines of code

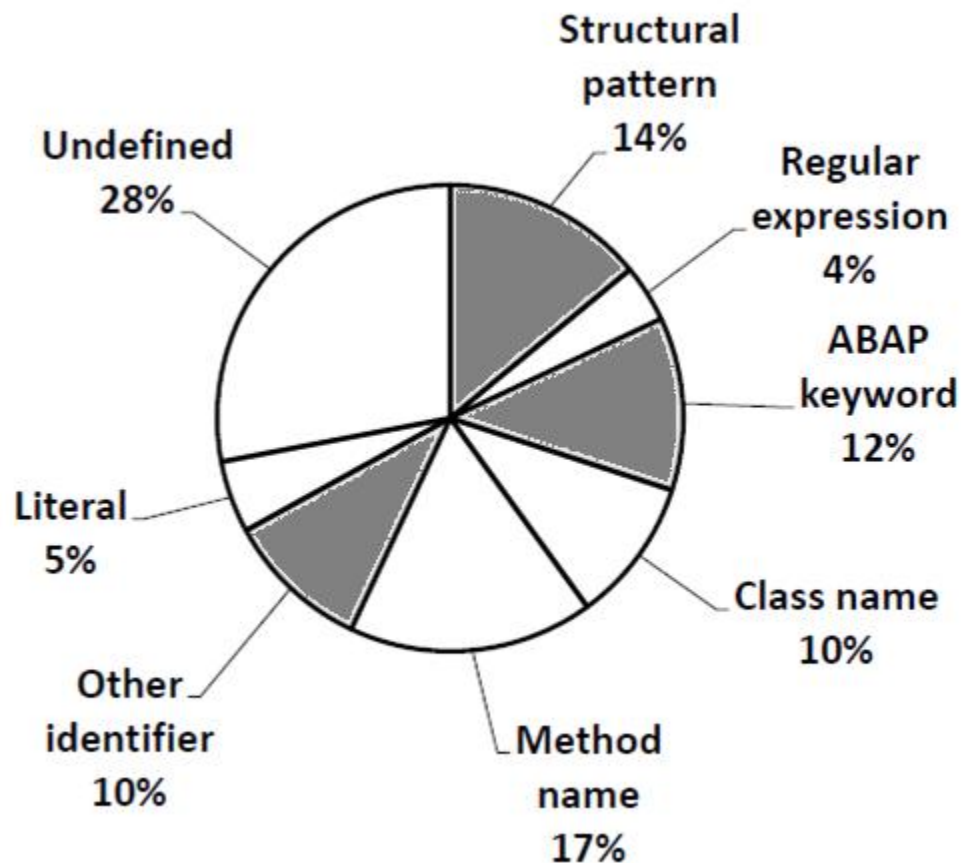
– at Google:

- 20+ Mio lines of code and 900k projects – released for public access

– on the Web:

- SourceForge.net – 295k open source projects
- Google Project Hosting – 250k open source projects

Categories of Developers' Queries



Example search targets:

1. compute statements that assign a value to a global variable
2. updating access to a database table
3. violation of some development guidelines
4. is a modification of (3) with an identifier added
5. all function invocations which use, e.g., `param_name`, as a formal parameter name
6. all occurrences of a system/global variable, e.g., `uname`, in an `ASSERT` statement
7. data declarations of a certain type
8. authorization checks

```

SELECT SINGLE * FROM bkpf INTO lwa_bkpf
WHERE bukers = bsik-bukers
AND belnr = bsik-belnr
AND gjahr = bsik-gjahr.

ktab-xbinr_alt = lwa_bkpf-xbinr_alt.

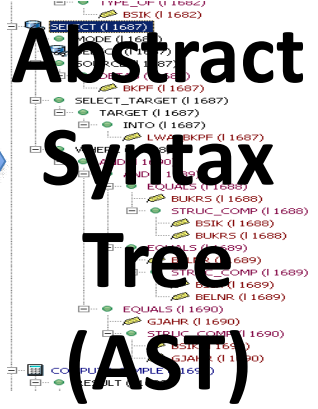
* Invoices and credit memos
IF bsik-zmsks IS INITIAL AND
bsik-zmsk IS INITIAL AND
* Note 512749 - EoI
bsik-rebzj IS INITIAL
* Note 512749 - EoI
* Note 512749 - EoI
bsik-xzahl IS INITIAL AND
bsik-bschl IN buschl.
* Note 512749 - EoI
APPEND ktab TO evt_invoice_tab.
ENDIF.

* Downpayment clearings (part pointing)
IF bsik-zmsks IS INITIAL AND
bsik-zmsk IS INITIAL AND
bsik-rebzj IS INITIAL AND
bsik-rebzg IS INITIAL AND
bsik-rebz IS INITIAL AND
bsik-rebz IS INITIAL.
APPEND ktab TO evt_downpay_tab.
ENDIF.

* Downpayment clearings (part pointing

```

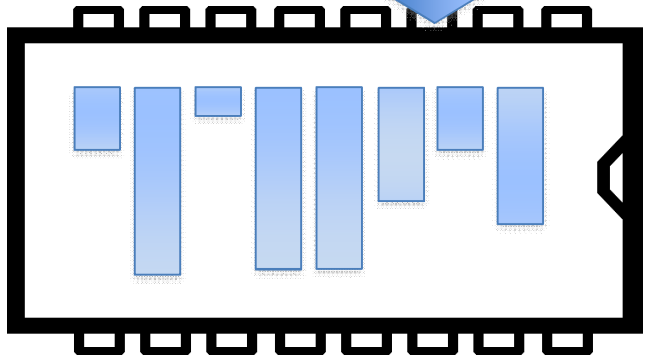
Source code



Abstract Syntax Tree (AST)



	Attribute	# of Rows	Cardinality
PK	DEV_OBJECT_NAME	458M	124k
	VERTEX_VALUE	458M	5,7M
	VERTEX_CATEGORY	458M	6
PK	PRE_ORDER	458M	5,3M
	POST_ORDER	458M	5,3M
	INCLUDE_NAME	458M	1,6M
	SOURCE_POSITION	458M	115k
FK	PARENT_PRE_ORDER	458M	3,4M



CALL FUNCTION ...
EXPORTING simplified = ...



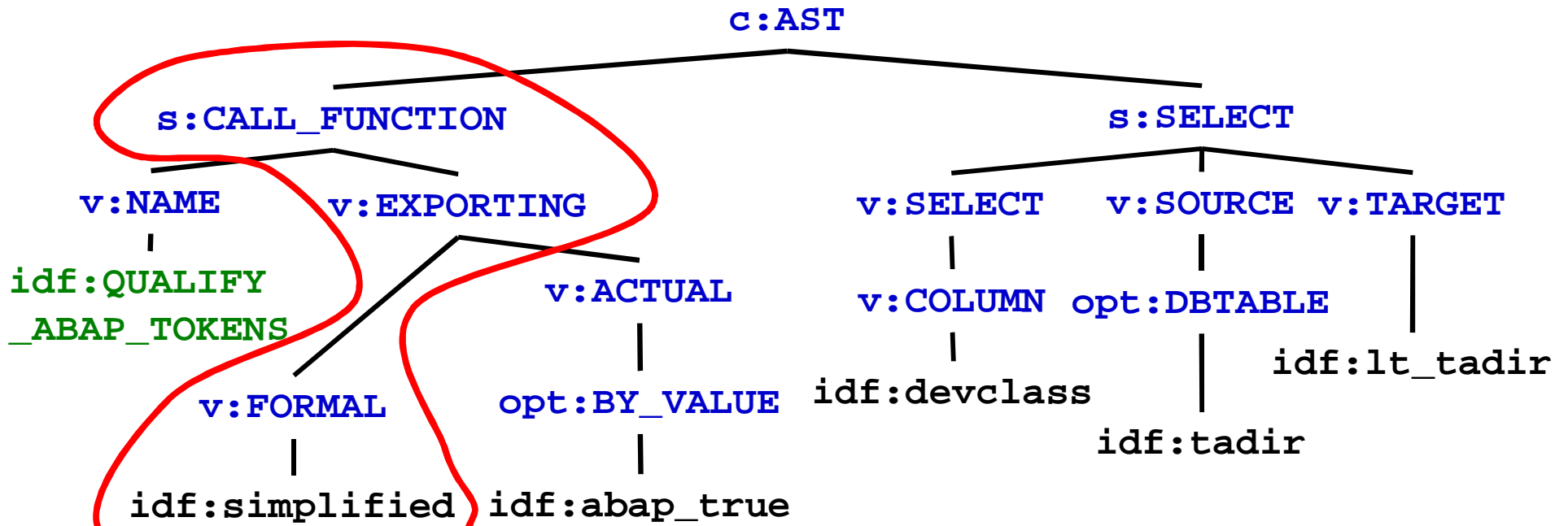
```

//s:CALL_FUNCTION
[./v:EXPORTING/v:FORMAL/idf:simplified]

```



```
CALL FUNCTION 'QUALIFY_ABAP_TOKENS'
  EXPORTING simplified = abap_true.
SELECT devclass FROM tadir INTO lt_tadir.
```

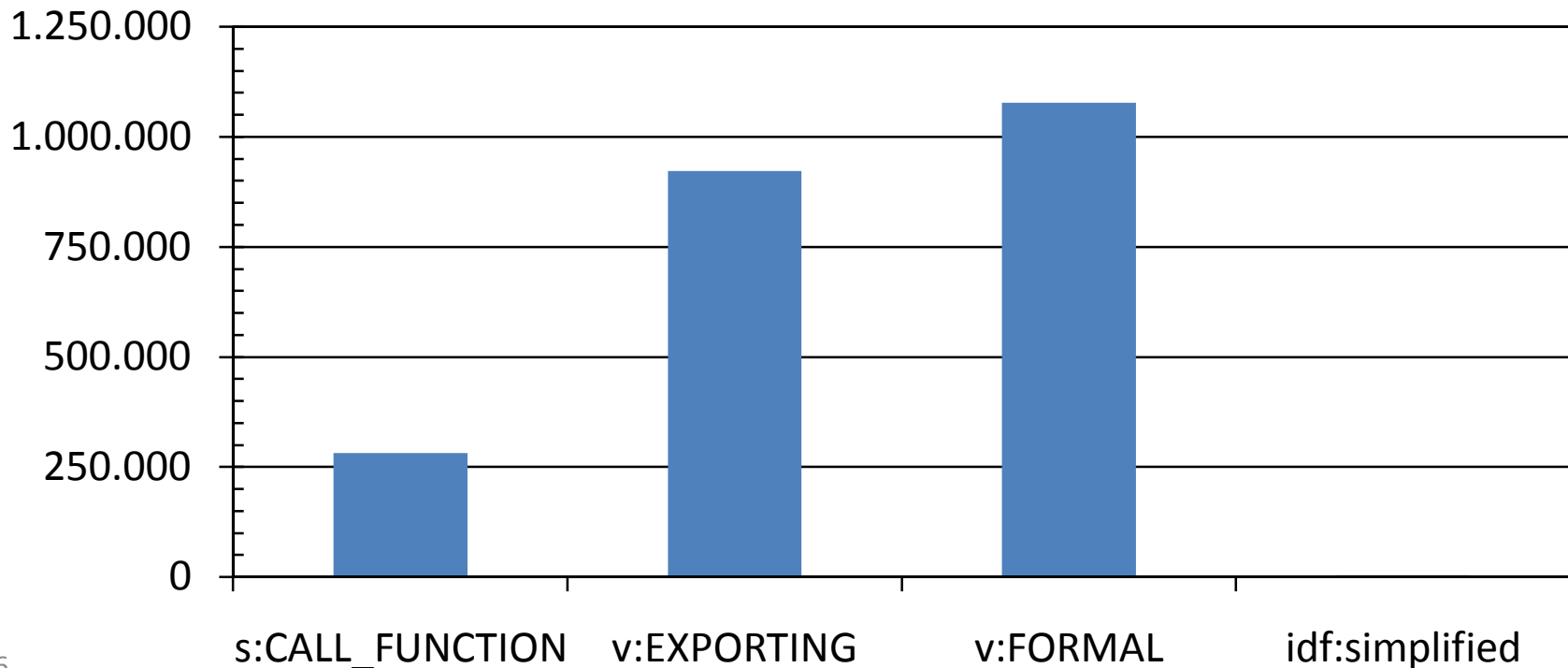


DEV_OBJECT_NAME	VERTEX_VALUE	VERTEX_CATEGORY	PRE_ORDER	POST_ORDER	...
CL_FREE_SEARCH	AST	compound	1	38	...
CL_FREE_SEARCH	CALL_FUNCTION	statement	2	19	...
CL_FREE_SEARCH	NAME	clause	3	6	...
CL_FREE_SEARCH	QUALIFY_ABAP_TOKENS	literal	4	5	...
CL_FREE_SEARCH	EXPORTING	clause	7	18	...
CL_FREE_SEARCH	FORMAL	clause	8	11	...
CL_FREE_SEARCH	simplified	identifier	9	10	...
...

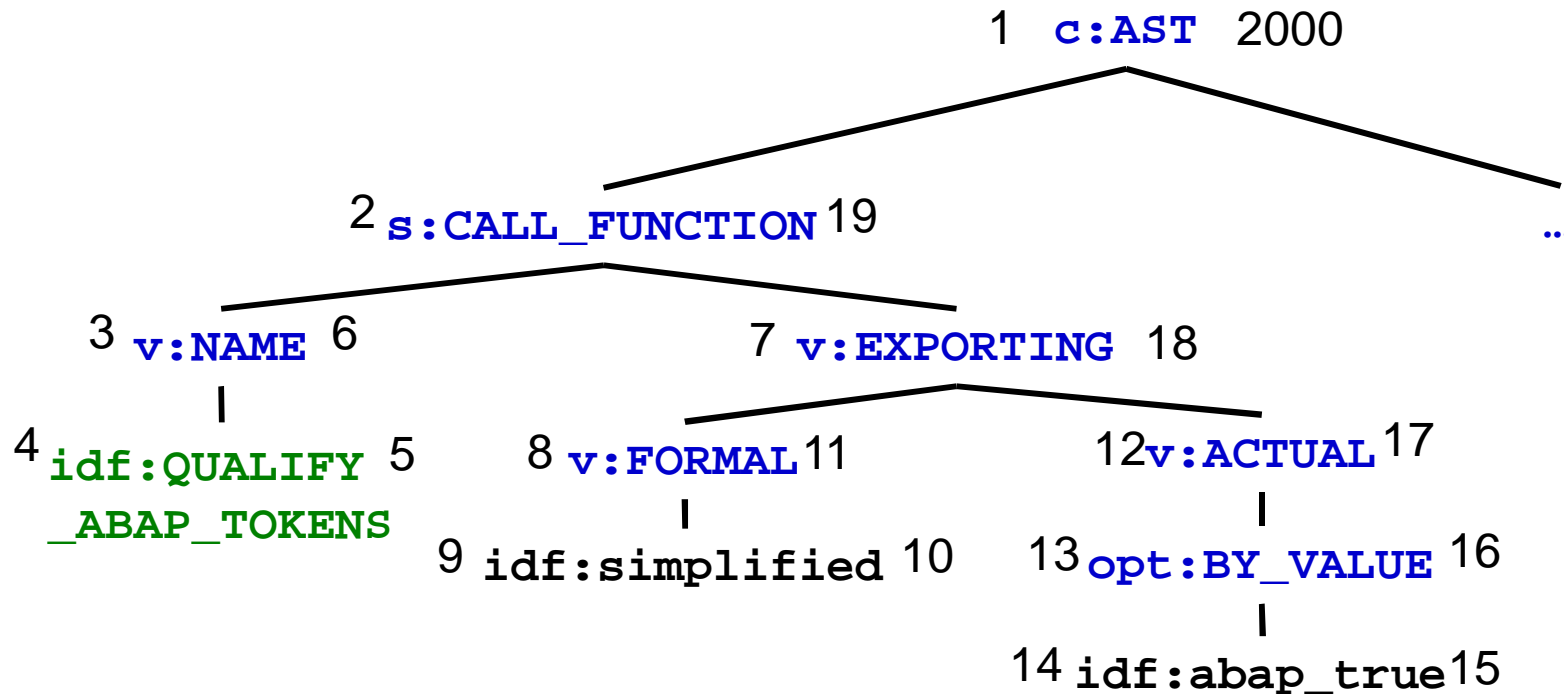
Query Execution

```
//s:CALL_FUNCTION[ ./v:EXPORTING/v:FORMAL/idf:simplified]
```

Table scan result size



Pre- Post-Order Numbering



Why In-Memory?

- **Since the number of AST vertex types is limited, the database can be well compressed**
 - **3,2GB of original source code in text files**
 - **80,0GB of AST data because of parsing**
 - **8,5GB compressed AST data**
- **Special indexing allows faster query evaluation**
 - **Pre-post-order numbering**
 - **Set processing instead of traversal of ASTs**
 - **Parallel data processing**

ABAP - BAP_100_c5080669_en [BAP, 100, C5080669, EN]/SEDI_EXT/CL_SEDI_ABS_NAV_METHOD - ABAP Development Tools

File Edit Source Navigate Search Project Run Window Help

Project Explorer

- BAP_100_c5080669_en [BAP, 100, C5080669, EN]

CL_SEDI_ABS_NAV_METH X VER04857 VER06510 »13

```
2137 IF sy-subrc <> 0.  
2138 LOOP AT <l_token>-statements INTO l_statement WHERE from  
2139 CALL FUNCTION 'RS_QUALIFY_ABAP_TOKENS_STR'  
2140 EXPORTING  
2141 statement_type = l_statement-type  
2142 index_from = l_statement-from  
143 index_to = l_statement-to  
144 simplified = abap_true  
2145 CHANGING  
2146 stokesx_tab = <l_token>-tokens  
2147 EXCEPTIONS  
2148 OTHERS = 1.  
2149 IF sy-subrc = 0  
2150 LOOP AT <l_token>-tokens INTO l_token FROM l_statement-  
2151 SEPARATED INTO TABLE <l_token>-qualified_tk_ta  
2152 ENDLLOOP.  
2153 EXIT.
```

DEMO

Search

ABAP Code Search | ABAP Search | File Search

Search can be performed.

Search Term : *

CALL FUNCTION 'RS_QUALIFY_ABAP_TOKENS_STR' EXPORTING

Restrict number of results to 100

Advanced Search

Scope

- ABAP Project: BAP_100_c5080669_en [BAP, 100, C5080669, EN]
- All ABAP Projects

Search in:

- Favorite Packages
- Local Objects
- All Packages

Customize... Search Cancel

Problems Search

ABAP Code Search for "CALL FUNCTION 'RS_QUALIFY_ABAP_TOKENS_STR' EXPORTING SIMPLIFIED =" 12
12 results (1,63 seconds)

Report Bug

- BAP_100_c5080669_en [BAP, 100, C5080669, EN]
- SEDI_EXT
 - CL_SEDI_ABS_NAV_METHOD
 - CL_SEDI_ABS_NAV_METHOD=====CM00G, Line Number: 13
 - CL_SEDI_PATTERN_GENERATOR
 - CL_SEDI_PATTERN_GENERATOR=====CM00D, Line Number: 56
- SEDI
- SVER_BC_ABA_LA_ACI_EAP
- SVER_BC_ABA_LA_VERIS

Read-Only Smart Insert 2139 : 49 BAP, 100, Id...C5080669, EN

Other Applications

- **Database table usage analysis**
- **Clone detection (visualisation)**
- **Trace analysis**

Database Table Usage Analysis

- **Goal: reengineering of the database layer**
- **For each table identify its usage patterns in source code (footprint):**
 - **In which packages/development objects is used**
 - **Which columns are used (projected, selected, joined, etc.)**
 - **What are technical settings of the table (buffering, sizing, total number of columns, etc.)**
 - **In which context is used (statements around the SQL statement)**
- **Find tables which follow a certain usage pattern, for example, large or wide tables, select statements in a loop, incorrect sizing, etc.**

