



SWAT

# Grammar Comparison Techniques

Vadim Zaytsev, SWAT, CWI  
Research Colloquium @ SERG  
2011

Questions?

# What this talk is not about

---

- ★ (Un)decidability of grammar equivalence
- ★ Model synchronisation
- ★ Tree diffing
- ★ Graph diffing

# Claims we all hear

---

- ★ “This compiler implements that language”
- ★ “This appendix contains an [*insert parsing technique here*] optimised grammar of the language”
- ★ “This grammarware produces data suitable to use with that grammarware”
- ★ “These are 100 implementations of one language”
- ★ “This language is a subset/superset of that language”
- ★ “This version of a compiler is backward compatible”

# Grammar differences

---

- ★ intended vs. accidental
- ★ result of grammar adaptation
- ★ result of grammar evolution
- ★ idiosyncrasies thanks to metanotation
- ★ idiosyncrasies thanks to parsing technology
- ★ presentation and understandability
- ★ misspelling
- ★ ...etc

# Grammar (non)equivalence

---

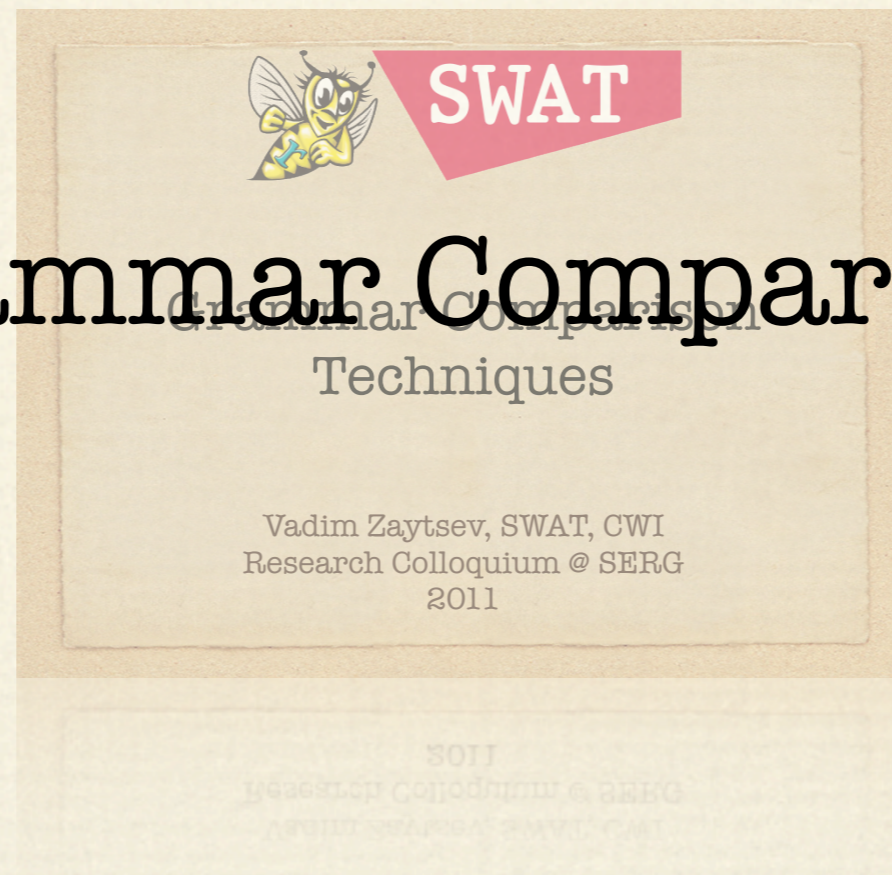
- ★ Undecidable
- ★ Non-existent
- ★ Can we cheat?
- ★ Approach (non)equivalence
- ★ Establish and maintain grammar relationships
- ★ Consistence throughout language incarnations!

# Conventions

---

## Language Comparison

### Grammar Comparison



# Assumed infrastructure

---

- ★ Grammar *format* free from idiosyncrasies
- ★ (Parse) *tree format* containing productions & literals
- ★ Grammar *extraction* for notation mapping:
  - ◆ Abstraction from predicates, sem.actions, rew.rules
- ★ Grammar *comparison* for spotting grammar differences
  - ◆ Nominal differences; structural differences
- ★ Grammar *transformation* suite:
  - ◆ Refactoring; extension / restriction; revision
- ★ Grammar *export* for connecting to different toolsets



# Straightforward comparison

---

★ Equivalence as equality

★ Nominal differences

◆  $A ::= X Y Z;$                        $B ::= X Y Z;$

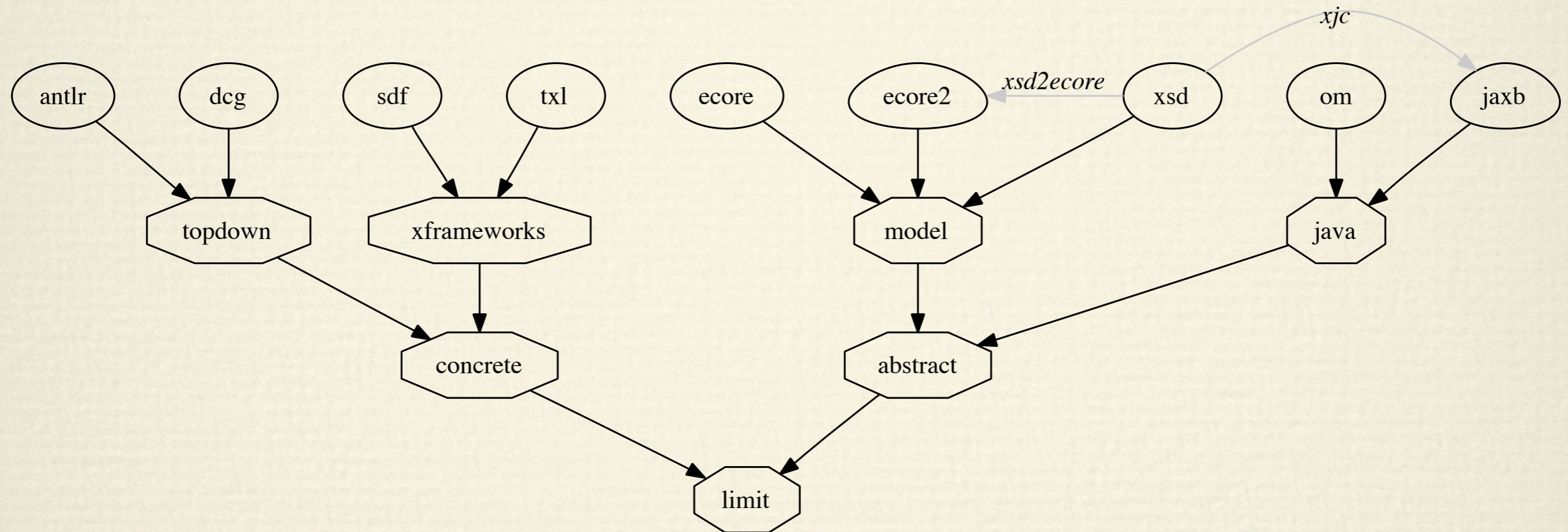
★ Structural differences

◆  $A ::= X Y Z;$                        $A ::= X \quad Z;$

★ Deliberately limited comparator is useful

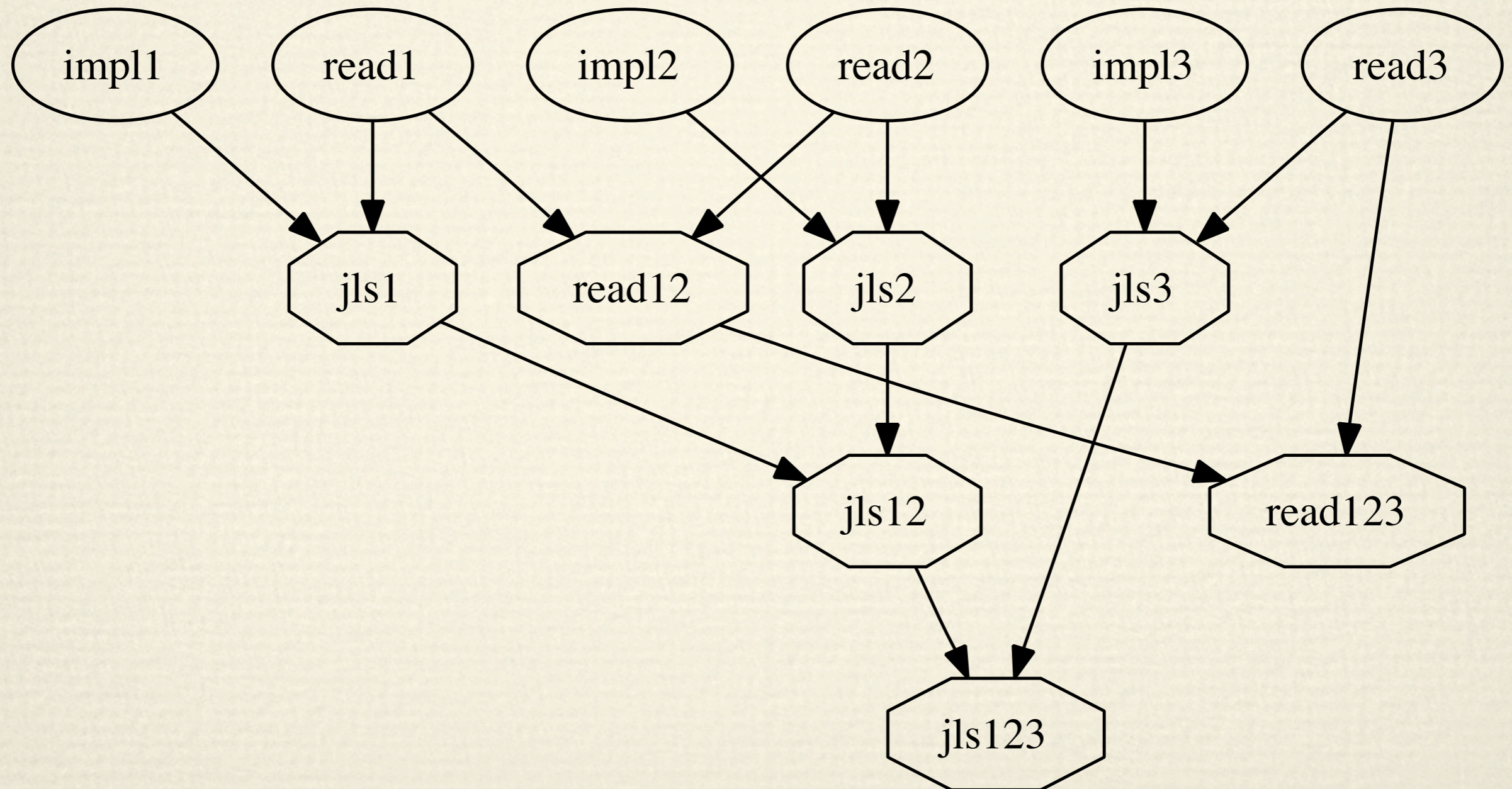
# Complicated scenario

Different implementations of the same language  
(parsers, data models, etc.)

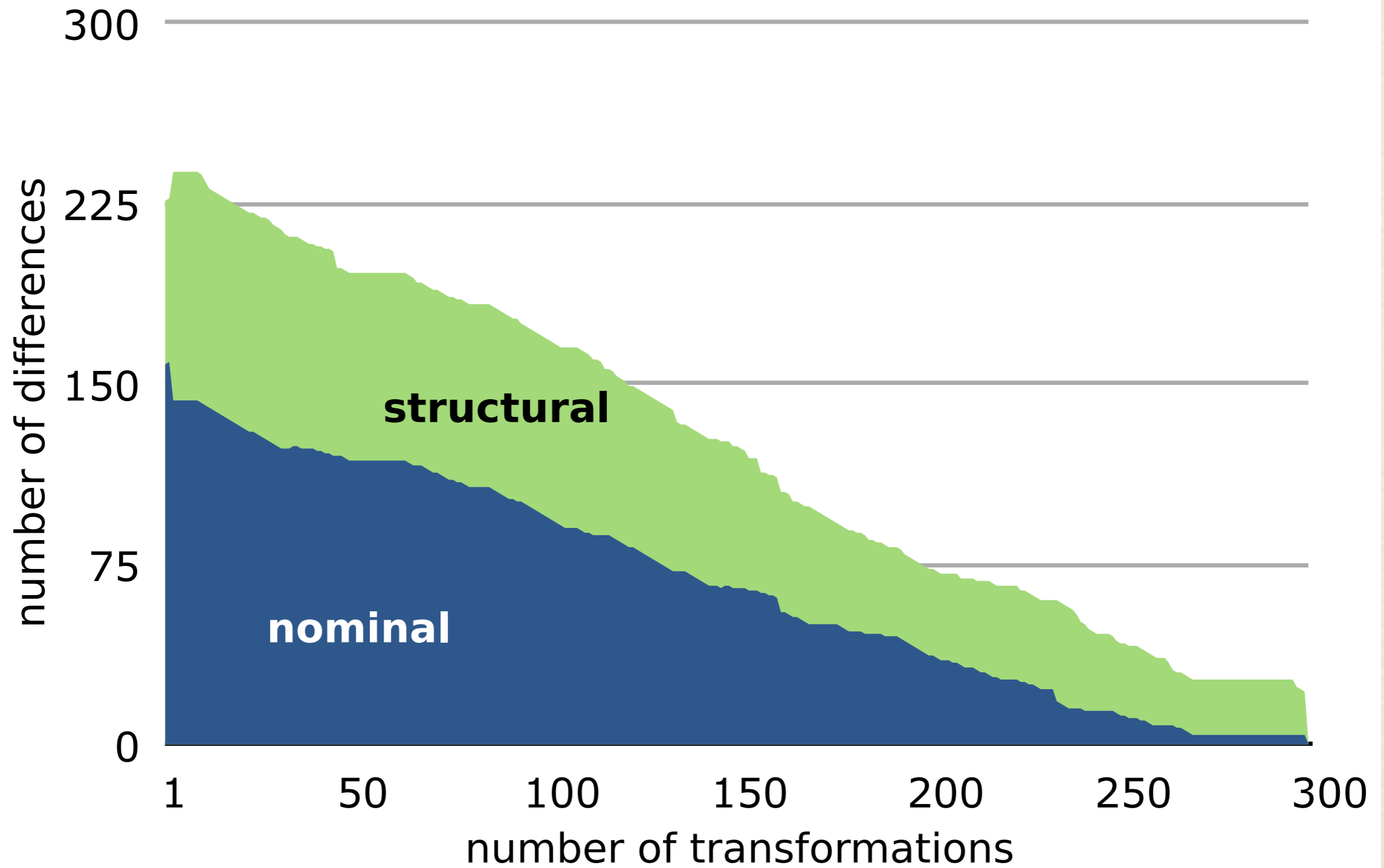


# Alternative scenario

Different versions of a language as documented by specifications



# Transform until equal



# Grammar refactoring example

---

BGF (*read2*)

ClassBodyDeclarations:  
  ClassBodyDeclaration

ClassBodyDeclarations:  
  ClassBodyDeclarations ClassBodyDeclaration

ClassBody:  
  "{" ClassBodyDeclarations? "}"

ClassBody:

"{" ClassBodyDeclaration \* "}"



XBGF (*grammar refactoring*)

**deyaccify**(ClassBodyDeclarations);

**inline**(ClassBodyDeclarations);

**massage**(

  ClassBodyDeclaration+? ,

  ClassBodyDeclaration \* );

# Grammar extension example

---

BGF (*read2*)

ClassModifier:

"public"  
"protected"  
"private"  
"abstract"  
"static"  
"final"  
"strictfp"

FieldModifier:

"public"  
"protected"  
"private"  
"static"  
"final"  
"transient"  
"volatile"

MethodModifier:

"public"  
"protected"  
"private"  
"abstract"  
"static"  
"final"  
"synchronized"  
"native"  
"strictfp"

XBGF (grammar optimisation)

```
unite(InterfaceModifier, Modifier);  
unite(ConstructorModifier, Modifier);  
unite(MethodModifier, Modifier);  
unite(FieldModifier, Modifier);
```

... ..

# Grammar revision example

---

BGF (*impl2, impl3*)

Expression2:

Expression3 Expression2Rest ?

Expression2Rest:

( Infixop Expression3 )\*

Expression2Rest:

~~Expression3~~ "instanceof" Type

XBGF (*grammar correction*)

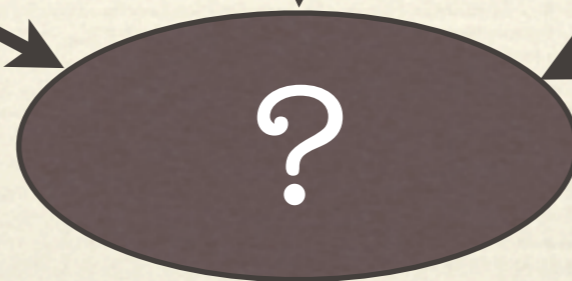
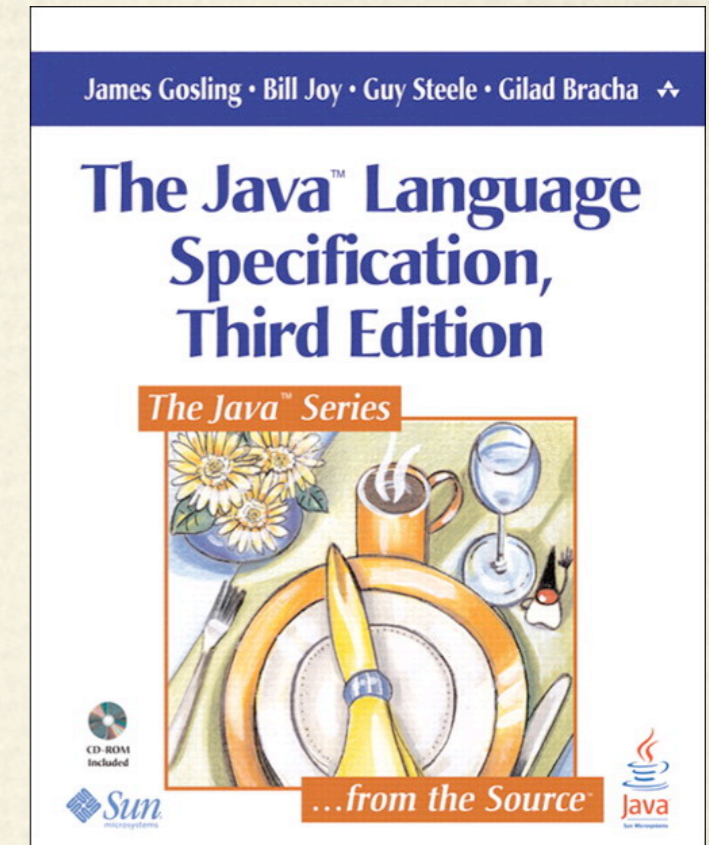
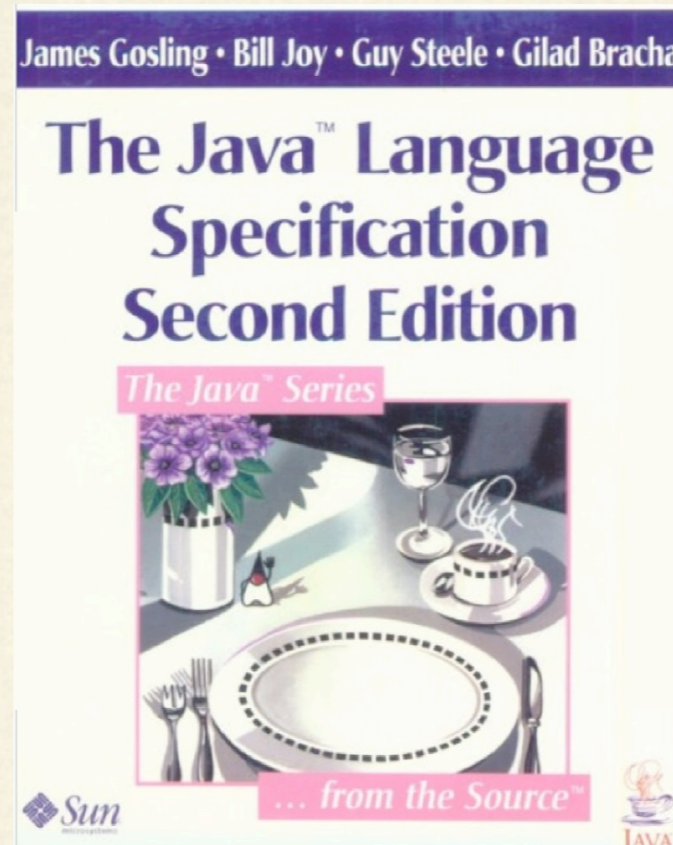
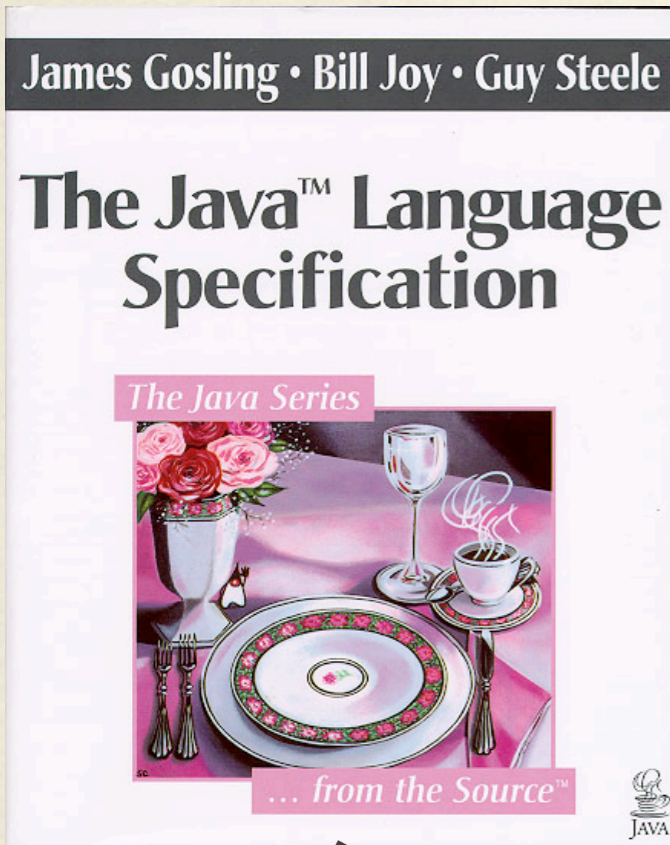
**project**(

Expression2Rest:

< Expression3 > "instanceof" Type

);

# Case study: JLS





# JLS convergence results

	<b>jls1</b>	<b>jls12</b>	<b>jls123</b>	<b>jls2</b>	<b>jls3</b>	<b>read12</b>	<b>read123</b>	<b>Total</b>
Number of lines	682	5114	2847	6774	10721	1639	3082	30859
Number of transformations	67	290	111	387	544	77	135	1611
○ Semantics-preserving (§4.2.2)	45	231	80	275	381	31	78	1121
○ Semantics-increasing/-decreasing	22	58	31	102	150	39	53	455
○ Semantics-revising	—	1	—	10	13	7	4	35
Preparation phase (§4.2.1)	1	—	—	15	24	11	14	65
○ Known bugs	—	—	—	1	11	—	4	16
○ Post-extraction	—	—	—	7	8	7	5	27
○ Initial correction	1	—	—	7	5	4	5	22
Resolution phase	21	59	31	97	139	35	43	425
○ Extension (§4.2.3)	—	17	26	—	—	31	38	112
○ Relaxation (§4.2.4)	18	39	5	75	112	—	2	251
○ Correction (§4.2.5)	3	3	—	22	27	4	3	62

Convergence reveals relationships

	<b>jls1</b>	<b>jls12</b>	<b>jls123</b>	<b>jls2</b>	<b>jls3</b>	<b>read12</b>	<b>read123</b>	<b>Total</b>
○ <i>rename</i>	9	4	2	9	10	—	2	36
○ <i>reroot</i>	2	—	—	2	2	2	1	9
○ <i>unfold</i>	1	10	8	11	13	2	3	48
○ <i>fold</i>	4	11	4	11	13	2	5	50
○ <i>inline</i>	3	67	8	71	100	—	1	250
○ <i>extract</i>	—	17	5	18	30	—	5	75
○ <i>chain</i>	1	—	2	—	—	1	4	8
○ <i>massage</i>	2	13	—	15	32	5	3	70
○ <i>distribute</i>	3	4	2	3	6	—	—	18
○ <i>factor</i>	1	7	3	5	24	3	1	44
○ <i>deyaccify</i>	2	20	—	25	33	4	3	87
○ <i>yaccify</i>	—	—	—	—	1	—	1	2
○ <i>eliminate</i>	1	8	1	14	22	—	—	46
○ <i>introduce</i>	—	1	30	4	13	3	34	85
○ <i>import</i>	—	—	2	—	—	—	1	3
○ <i>vertical</i>	5	7	7	8	22	5	8	62
○ <i>horizontal</i>	4	19	5	17	31	4	4	84
○ <i>add</i>	1	14	13	7	20	28	20	103
○ <i>appear</i>	—	8	11	8	25	2	17	71
○ <i>widen</i>	1	3	—	1	8	1	3	17
○ <i>upgrade</i>	—	8	—	14	20	2	2	46
○ <i>unite</i>	18	2	—	18	21	5	4	68
○ <i>remove</i>	—	10	1	11	18	—	1	41
○ <i>disappear</i>	—	7	4	11	11	—	—	33
○ <i>narrow</i>	—	—	1	—	4	—	—	5
○ <i>downgrade</i>	—	2	—	8	3	—	—	13
○ <i>define</i>	—	6	—	4	9	1	6	26
○ <i>undefine</i>	—	3	—	5	3	—	—	11
○ <i>redefine</i>	—	3	—	8	7	6	2	26
○ <i>inject</i>	—	—	—	2	4	—	1	7
○ <i>project</i>	—	1	—	1	2	—	—	4
○ <i>replace</i>	3	1	2	3	6	1	1	17
○ <i>unlabel</i>	—	—	—	—	—	—	2	2

# Resources

---

- ★ R. Lämmel, V. Zaytsev, *An Introduction to Grammar Convergence*, iFM 2009, LNCS 5423.
- ★ R. Lämmel, V. Zaytsev, *Reverse Engineering Grammar Relationships*, WSR 2010.
- ★ V. Zaytsev, *Language Convergence Infrastructure*, GTTSE 2009, LNCS 6491.
- ★ R. Lämmel, V. Zaytsev, *Recovering Grammar Relationships for the JLS*, SCAM 2009, SQJ 19:2, arXiv abs/1008.4188.

# Grammar-based testing

- ★ Asymmetric comparison:

- ◆ Reference grammar vs. parser under test

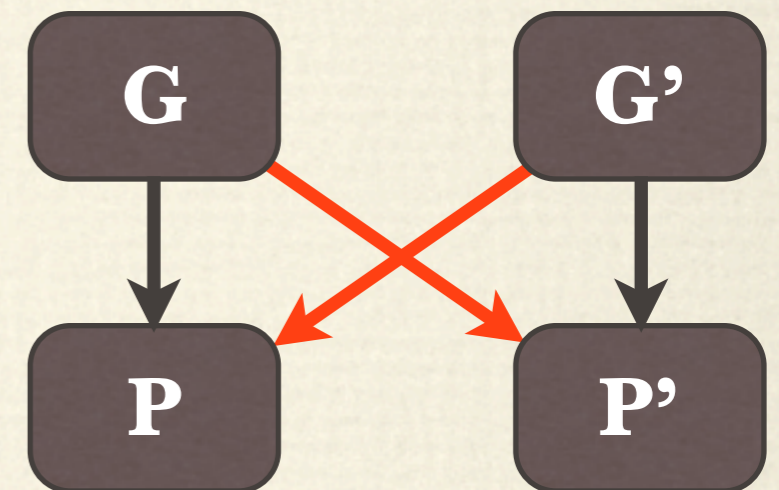
- ★ Symmetric comparison:

- ◆ Differential testing

- ★ Systematic test data generation

- ◆ Controlled combinatorial coverage

- ★ Larger sets of smaller test data items



# Test data generation (1/4)

$grammar(Ps)$   
 $\Leftarrow maplist(prod, Ps).$

$prod(p(L, N, X))$   
 $\Leftarrow mapopt(atom, L), atom(N), expr(X).$

$expr(true).$   
 $expr(t(T)) \Leftarrow atom(T).$   
 $expr(n(N)) \Leftarrow atom(N).$   
 $expr(','(Xs)) \Leftarrow maplist(expr, Xs).$   
 $expr(';')(Xs) \Leftarrow maplist(expr, Xs).$   
 $expr('?'(X)) \Leftarrow expr(X).$   
 $expr('*')(X) \Leftarrow expr(X).$   
 $expr('+'(X)) \Leftarrow expr(X).$

$tree(true).$   
 $tree(t(T)) \Leftarrow atom(T).$   
 $tree(n(P, T)) \Leftarrow prod(P), tree(T).$   
 $tree(','(Ts)) \Leftarrow maplist(tree, Ts).$   
 $tree(';')(X, T) \Leftarrow expr(X), tree(T).$   
 $tree('?'(Ts)) \Leftarrow mapopt(tree, Ts).$   
 $tree('*')(Ts) \Leftarrow maplist(tree, Ts).$   
 $tree('+'(Ts)) \Leftarrow maplist1(tree, Ts).$

# Test data generation (2/4)

$mark(C, p(L, N, X1), p(L, N, X2)) \Leftarrow$   
 $mark(C, X1, X2).$

Marked productions are essentially marked expressions.

$mark(uc, n(N), \{n(N)\}).$   
 $mark(bc, ';' (Xs), \{';' (Xs)\}).$   
 $mark(bc, '?' (X), \{'?' (X)\}).$   
 $mark(bc, '*' (X), \{'*' (X)\}).$   
 $mark(bc, '+' (X), \{'+' (X)\}).$

A nonterminal occurrence provides a focus for unfolding coverage. The EBNF forms ';', '?', '\*', '+' provide foci for branch coverage.

$mark(C, '?' (X1), '?' (X2)) \Leftarrow$   
 $mark(C, X1, X2).$   
 $mark(C, '*' (X1), '*' (X2)) \Leftarrow$   
 $mark(C, X1, X2).$   
 $mark(C, '+' (X1), '+' (X2)) \Leftarrow$   
 $mark(C, X1, X2).$

Foci for BC and UC may also be found by recursing into subexpressions.

$mark(C, ',' (Xs1), ',' (Xs2)) \Leftarrow$   
 $append(Xs1a, [X1|Xs1b], Xs1),$   
 $append(Xs1a, [X2|Xs1b], Xs2),$   
 $mark(C, X1, X2).$

Sequences and choices combine multiple expressions, and foci are found by considering one subexpression at the time.

$mark(C, ';' (Xs1), ';' (Xs2)) \Leftarrow$   
 $append(Xs1a, [X1|Xs1b], Xs1),$   
 $append(Xs1a, [X2|Xs1b], Xs2),$   
 $mark(C, X1, X2).$

# Coverage criteria

---

- ★ **Trivial** coverage: if the test data set is not empty.
- ★ **Nonterminal** coverage: if each nonterminal is exercised at least once.
- ★ **Production** coverage: if each production in the grammar is exercised at least once.
- ★ **Branch** coverage: each branch of  $? | * +$
- ★ **Unfolding** coverage: each production of each right hand side nonterminal occurrence
- ★ **Context-dependent branch coverage!**

# Test data generation (3/4)

$\text{vary}(G, \{n(N)\}, n(P, T)) \Leftarrow$   
 $\text{def}(G, N, Ps),$   
 $\text{member}(P, Ps),$   
 $P = p(-, -, X),$   
 $\text{complete}(G, X, T).$

A nonterminal occurrence in focus is varied so that all productions are exercised. (The complete spec also deals with chain productions and top-level choices in a manner that increases variation in a reasonable sense.)

$\text{vary}(G, \{';'(Xs)\}, ';'(X, T)) \Leftarrow$   
 $\text{member}(X, Xs),$   
 $\text{complete}(G, X, T).$

A choice in focus is varied so that all branches are exercised.

$\text{vary}(-, \{ '?'(-) \}, '?'([ ])).$

$\text{vary}(G, \{ '?'(X) \}, '?'([T])) \Leftarrow$   
 $\text{complete}(G, X, T).$

An optional expression and a '\*' repetition in focus are varied so that the cases for no tree and one tree are exercised. A '+' repetition is varied so that the cases for sequences of length 1 and 2 are exercised.

$\text{vary}(-, \{ '*'(-) \}, '*'([ ])).$

$\text{vary}(G, \{ '*'(X) \}, '*'([T])) \Leftarrow$   
 $\text{complete}(G, X, T).$

$\text{vary}(G, \{ '+'(X) \}, '+'([T])) \Leftarrow$   
 $\text{complete}(G, X, T).$

We omit all clauses for recursing into compound expressions; they mimic shortest completion but they are directed in a way that they reach the focus.

$\text{vary}(G, \{ '+'(X) \}, '+'([T1, T2])) \Leftarrow$   
 $\text{complete}(G, X, T1),$   
 $\text{complete}(G, X, T2).$



# Test data generation (4/4)

$tc(G, R, T)$

$\Leftarrow \text{def}(G, R, -), \text{complete}(G, n(R), T).$

$nc(G, R, T)$

$\Leftarrow \text{def}(G, R, -), \text{dist}(G, R, H, -), \text{hole}(G, n(R), H, T, V), \text{complete}(G, n(H), V).$

$pc(G, R, T)$

$\Leftarrow \text{def}(G, R, Ps), \text{member}(P, Ps), \text{complete}(G, P, T).$

$pc(G, R, T)$

$\Leftarrow \text{def}(G, R, -), \text{dist}(G, R, H, -), \text{hole}(G, n(R), H, T, V), \text{pc}(G, H, V).$

$bc(G, R, T)$

$\Leftarrow \text{cdbc}(bc, G, R, T).$

$uc(G, R, T)$

$\Leftarrow \text{cdbc}(uc, G, R, T).$

$\text{cdbc}(C, G, R, T)$

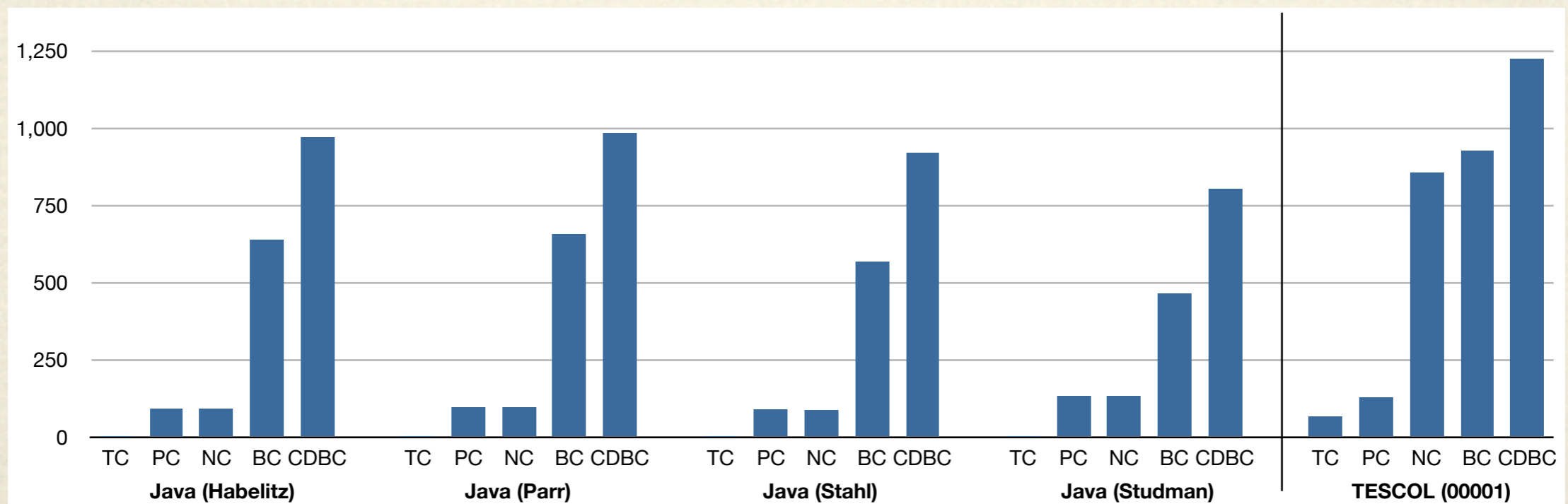
$\Leftarrow \text{def}(G, R, Ps), \text{member}(P, Ps), \text{mark}(C, P, F), \text{vary}(G, F, T).$

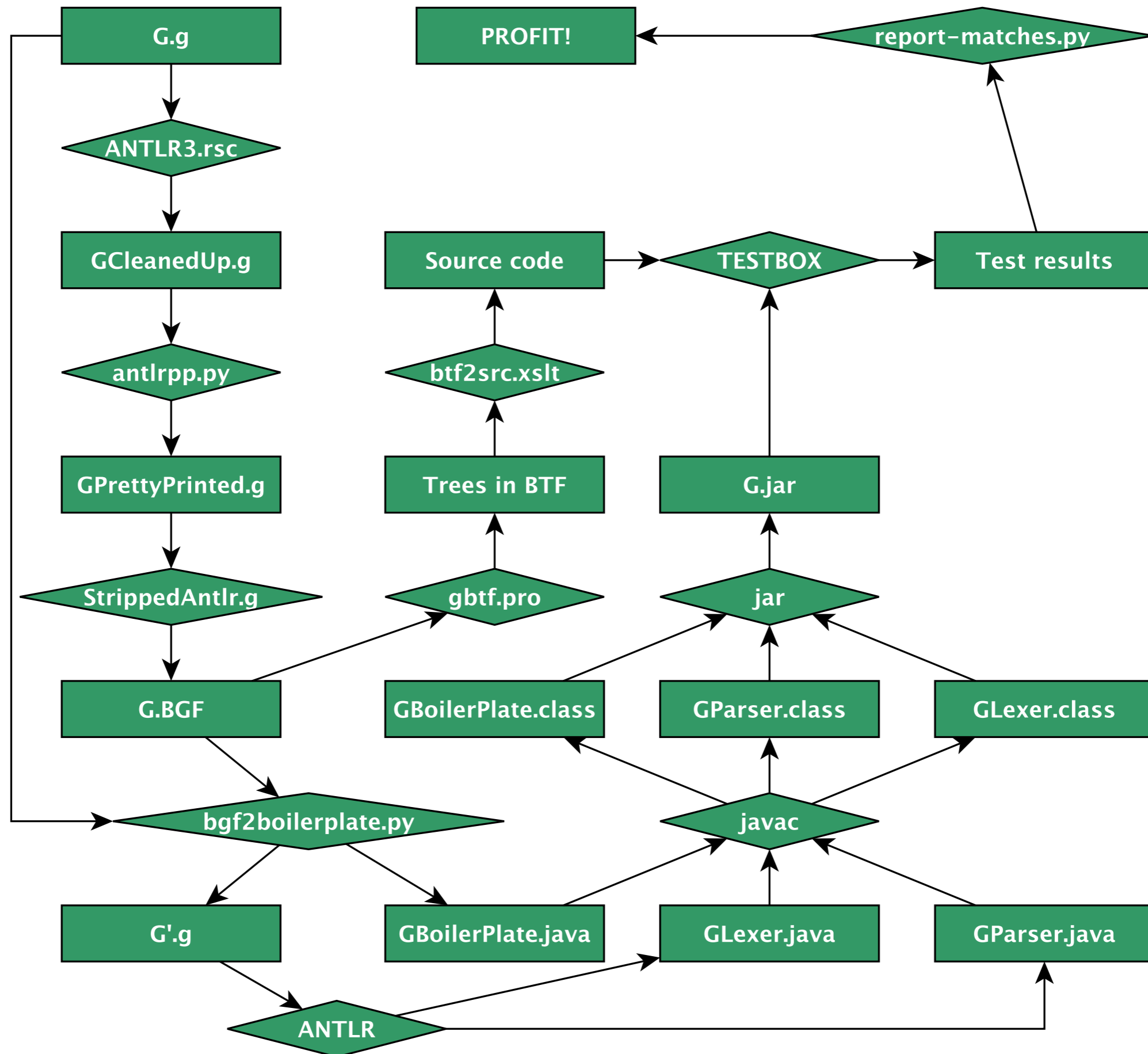
$\text{cdbc}(C, G, R, T)$

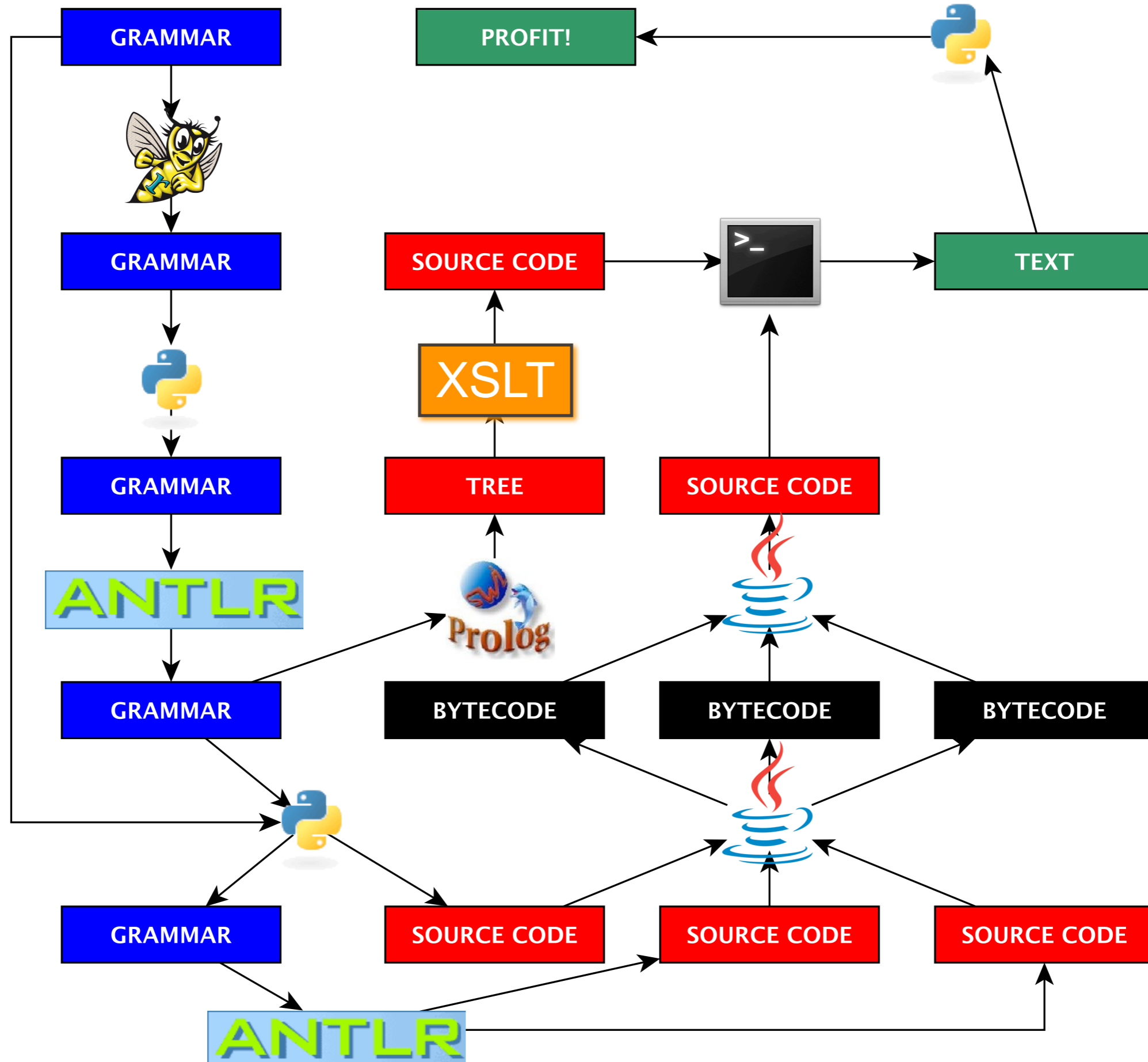
$\Leftarrow \text{def}(G, R, -), \text{dist}(G, R, H, -), \text{hole}(G, n(R), H, T, V), \text{cdbc}(C, G, H, V).$

# Grammar equivalence study: Java

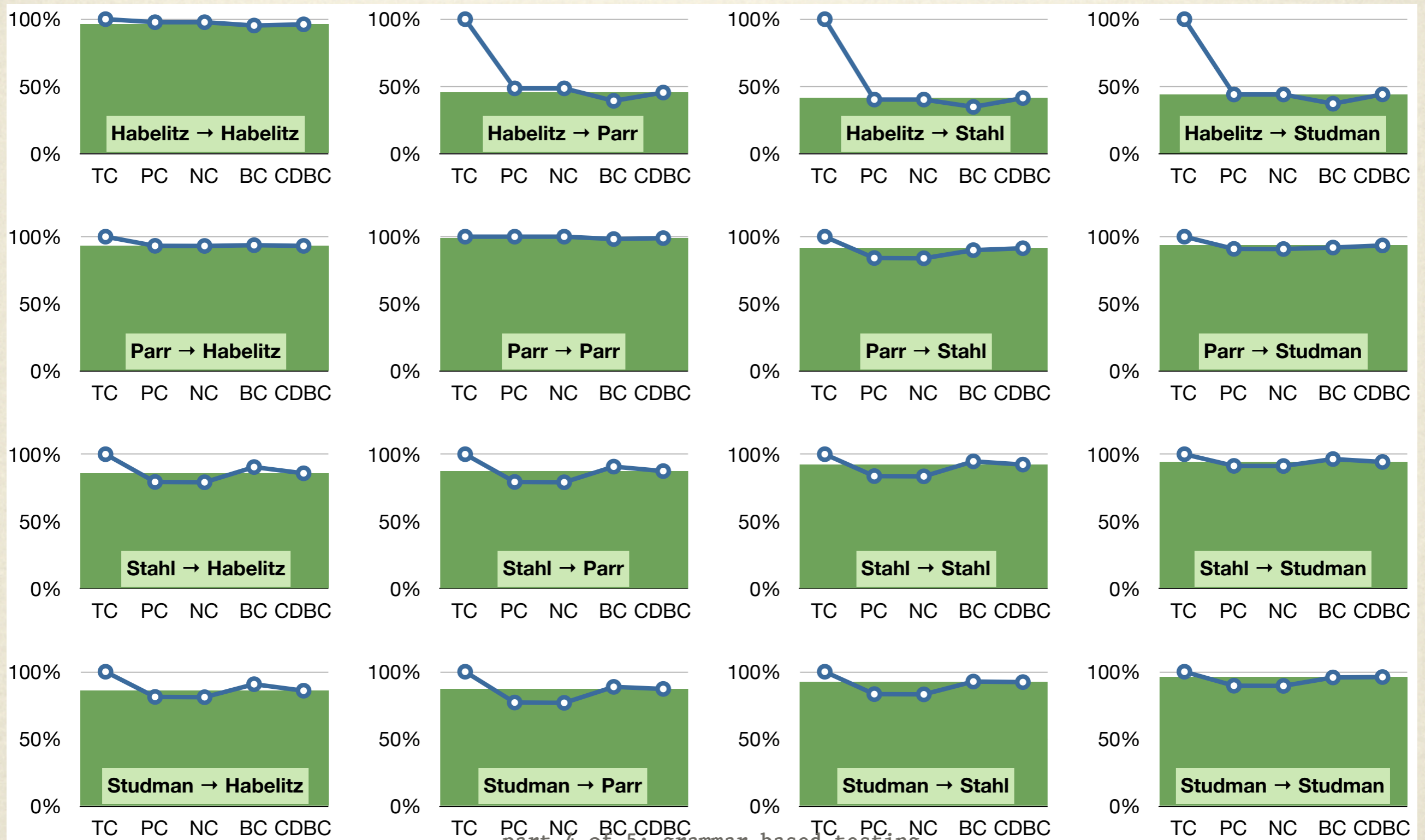
Codename	Tech	Author	year	PROD	VAR	TERM	...
<b>Habelitz</b>	ANTLR3	Dieter Habelitz	2008	397	226	166	...
<b>Parr</b>	ANTLR3	Terence Parr	2006	425	151	157	...
<b>Stahl</b>	ANTLR2	Michael Stahl	2004	262	155	167	...
<b>Studman</b>	ANTLR2	Michael Studman	2004	267	161	168	...







# Grammar equivalence study: Java



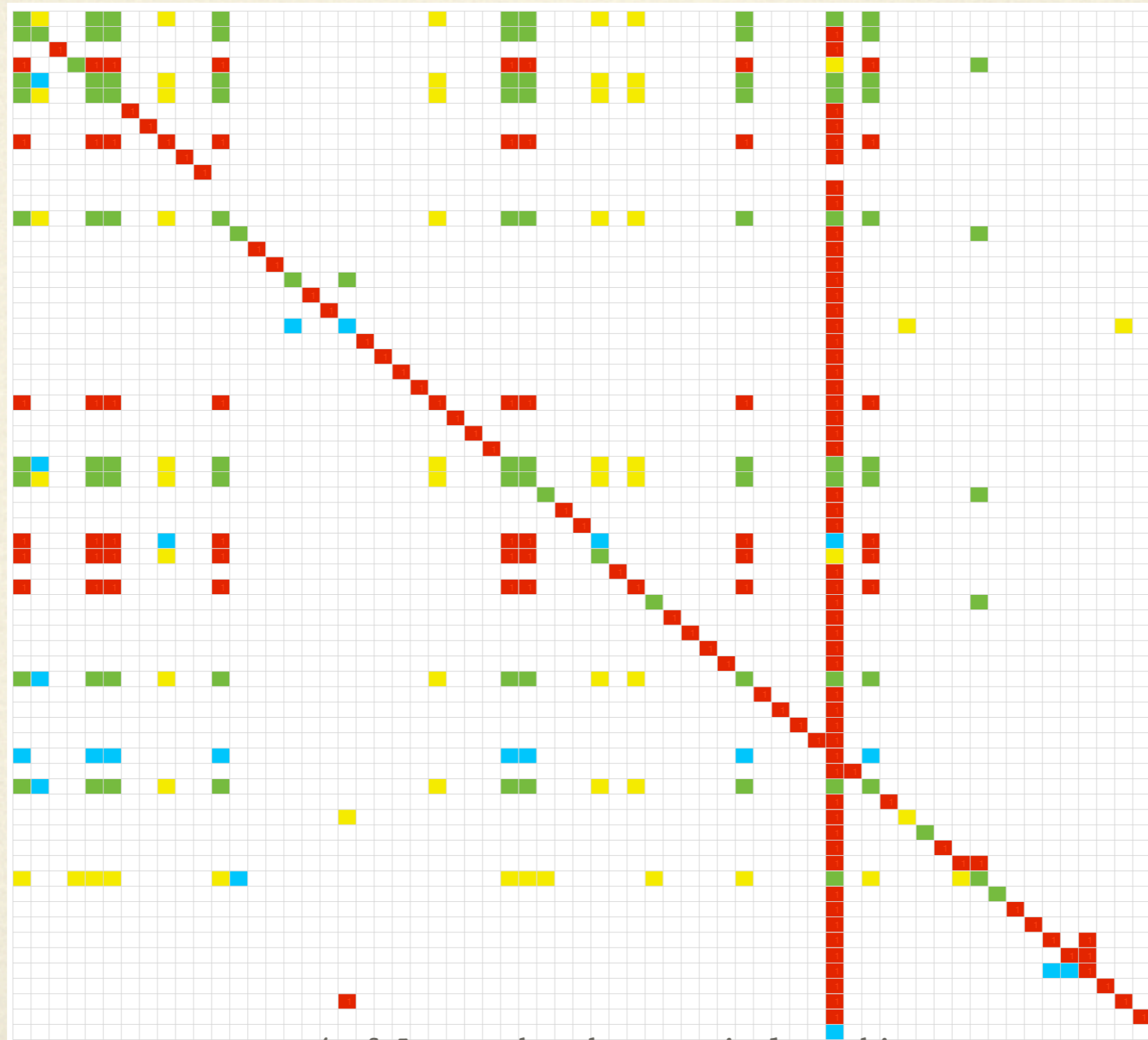
# Resources

---

- ★ B. Fischer, R. Lämmel, V. Zaytsev, *Comparison of Context-free Grammars Based on Parsing Generated Test Data*, SLE 2011, LNCS 6940.
- ★ ~~A. Davenport, B. Fischer, CC 2012 draft.~~

- <http://softlang.uni-koblenz.de/testmatch>
- <http://slps.sourceforge.net/testmatch>
- <http://grammarware.net/text/2011/testmatch.pdf>

# Nonterminal matching



part 4 of 5: test-based nonterminal matching

# Name matching study: TESCOL

.....

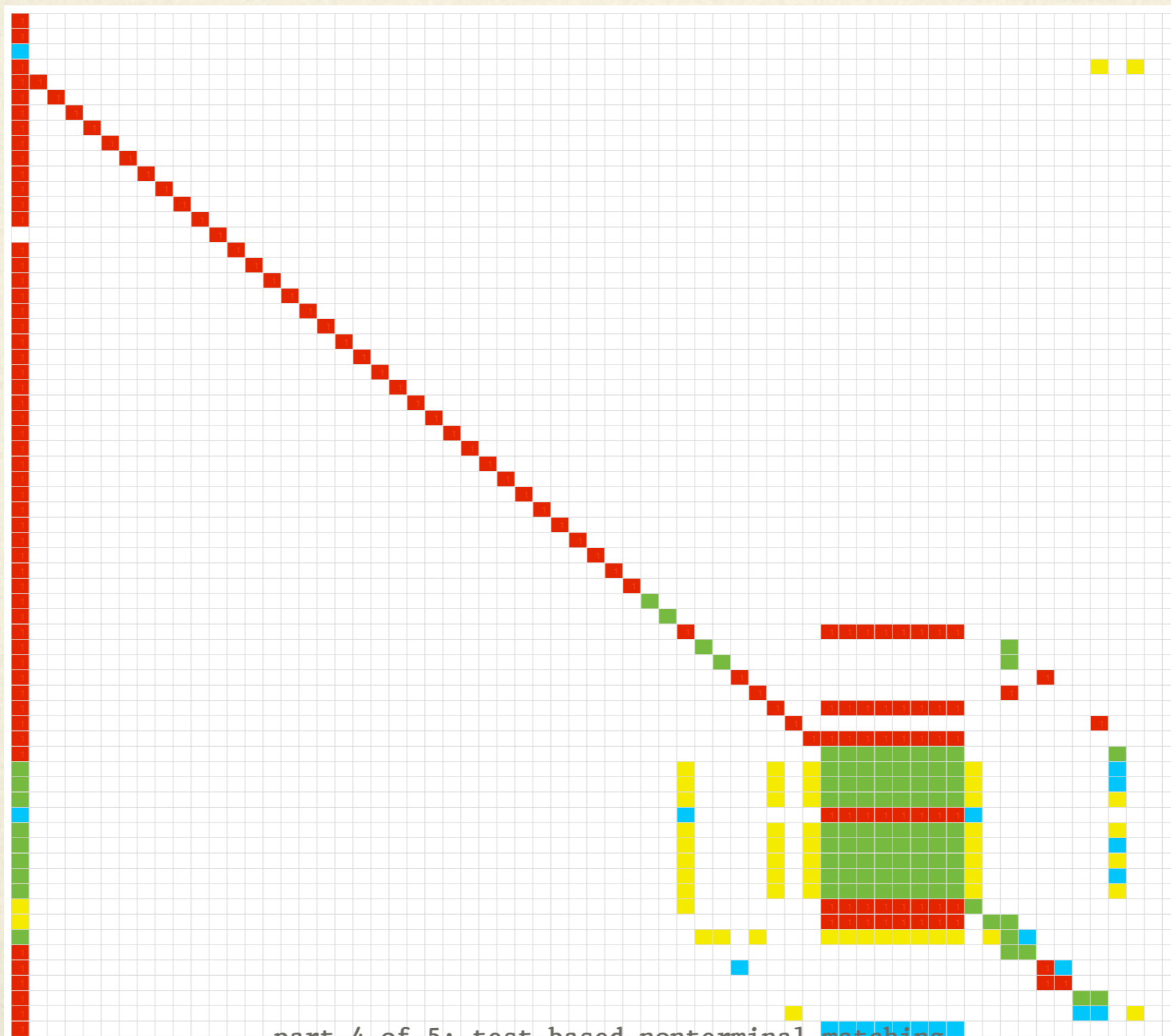
<b>Codename</b>	<b>Tech</b>	<b>Author</b>	<b>year</b>	<b>PROD</b>	<b>VAR</b>	<b>TERM</b>	<b>...</b>
<b>00000</b>	ANTLR3	[obfuscated]	2010	126	74	107	...
<b>00001</b>	ANTLR3	[obfuscated]	2010	79	67	107	...
<b>00010</b>	ANTLR3	[obfuscated]	2010	101	73	108	...
<b>00011</b>	ANTLR3	[obfuscated]	2010	84	63	107	...
<b>00100</b>	ANTLR3	[obfuscated]	2010	93	76	108	...
<b>00101</b>	ANTLR3	[obfuscated]	2010	94	76	107	...
<b>00110</b>	ANTLR3	[obfuscated]	2010	92	75	120	...
<b>00111</b>	ANTLR3	[obfuscated]	2010	84	71	108	...
<b>01000</b>	ANTLR3	[obfuscated]	2010	85	67	107	...
...	...	...	...	...	...	...	...



# Nonterminal matching



# Nonterminal matching



part 4 of 5: test-based nonterminal matching

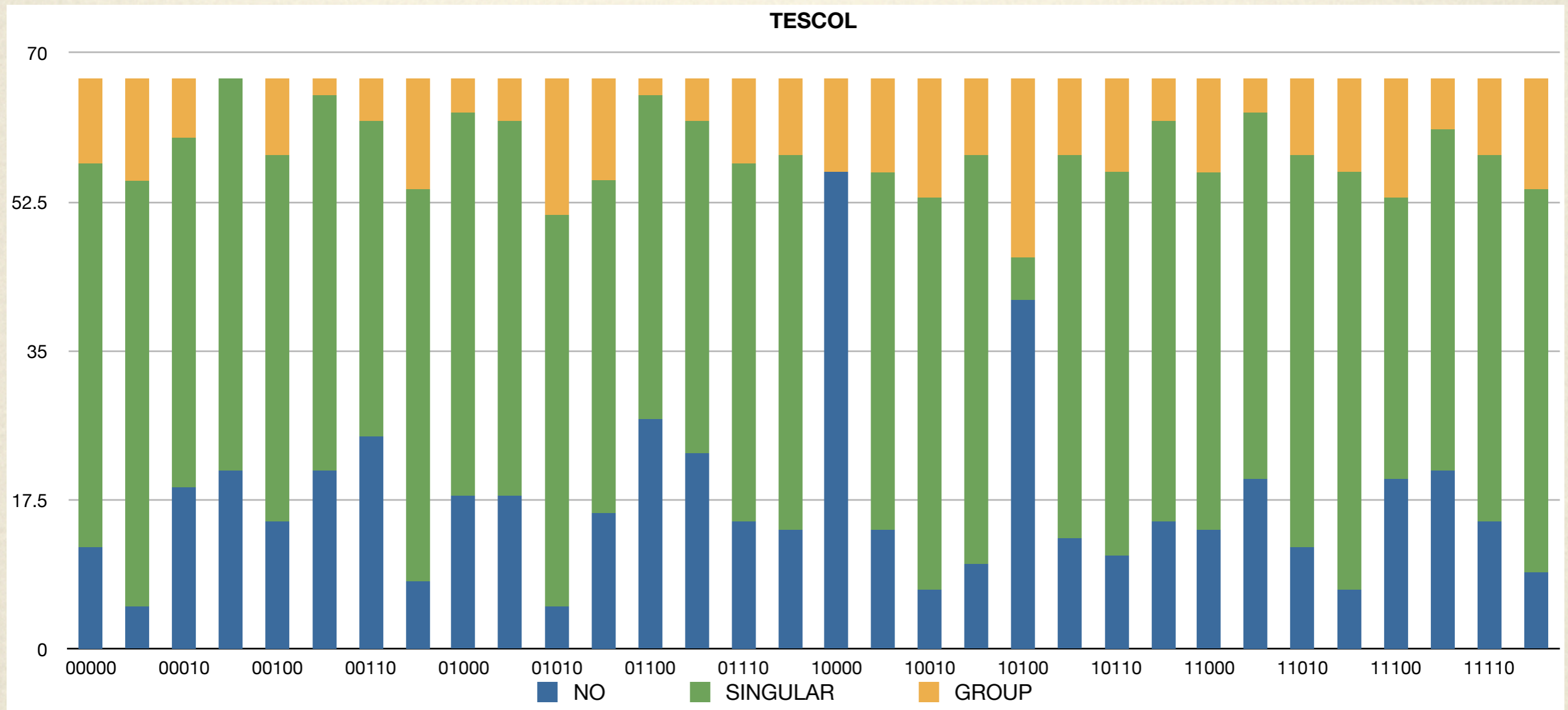
# Nonterminal matching



part 4 of 5: test-based nonterminal matching

# Nonterminal matching

.....



# Resources

---

★ B. Fischer, R. Lämmel, V. Zaytsev, *Comparison of Context-free Grammars Based on Parsing Generated Test Data*, SLE 2011, LNCS 6940.

★ ...

- <http://softlang.uni-koblenz.de/testmatch>
- <http://slps.sourceforge.net/testmatch>
- <http://slps.sourceforge.net/tank/#tescol>
- <http://grammarware.net/text/2011/testmatch.pdf>

# Roadmap

---

- ★ Motivation & methodology
- ★ Techniques not covered by the talk
- ★ Straightforward grammar comparison
- ★ Grammar convergence
- ★ Grammar-based testing
- ★ Test-based nonterminal matching

# Thank you!

- ★ <http://grammarware.net/slides/2011/comparison-serg.pdf>
- ★ <http://grammarware.net>
- ★ @grammarware
- ★ ...