

Generating Incremental Parsers

Christoph Höger – TU-Berlin

-
- „An LR(0) item (item for short) of a grammar G is a production of G with a dot at some position of the body.“
[Compilers – Principles, Techniques & Tools]
 - This may not be the best method to teach LR
 - PaGe (**P**arser **G**enerator) uses „Actions“:

$$A \rightarrow w \Rightarrow A \rightarrow w \textcircled{1}$$

-
- PaGe transforms the input Grammar
 - 1st Target: „Normalform“:

$$A \rightarrow t \underline{Z}_m$$

... (L(G) is kept unchanged)

$$Z_n \rightarrow \textcircled{1}$$

- 2nd Target: Erase original productions:

$$\begin{array}{ccc}
 A \rightarrow \underline{t} Z_m & & A \rightarrow \underline{t} Z_m \\
 \dots & \xrightarrow{\hspace{2cm}} & \dots \\
 B \rightarrow A Z_k & & B \rightarrow \underline{A} Z_k \\
 & & B \rightarrow \underline{t} Z_m Z_k
 \end{array}$$

- Note the Pseudo Terminal A
- A production can be deleted!

-
- 3rd Target: Remove Left recursion, leftfactor common start symbols
 - L(G) has not been changed yet!

$$Z_n \rightarrow \underline{A} Z_k$$

All rules start with (Pseudo-)

$$Z_n \rightarrow \underline{t} Z_m Z_k$$

Terminals (*shift*) or Actions (*reduce*)

$$Z_n \rightarrow \textcircled{1}$$

- This implies a LR Parse table

Incremental Parsing - What?

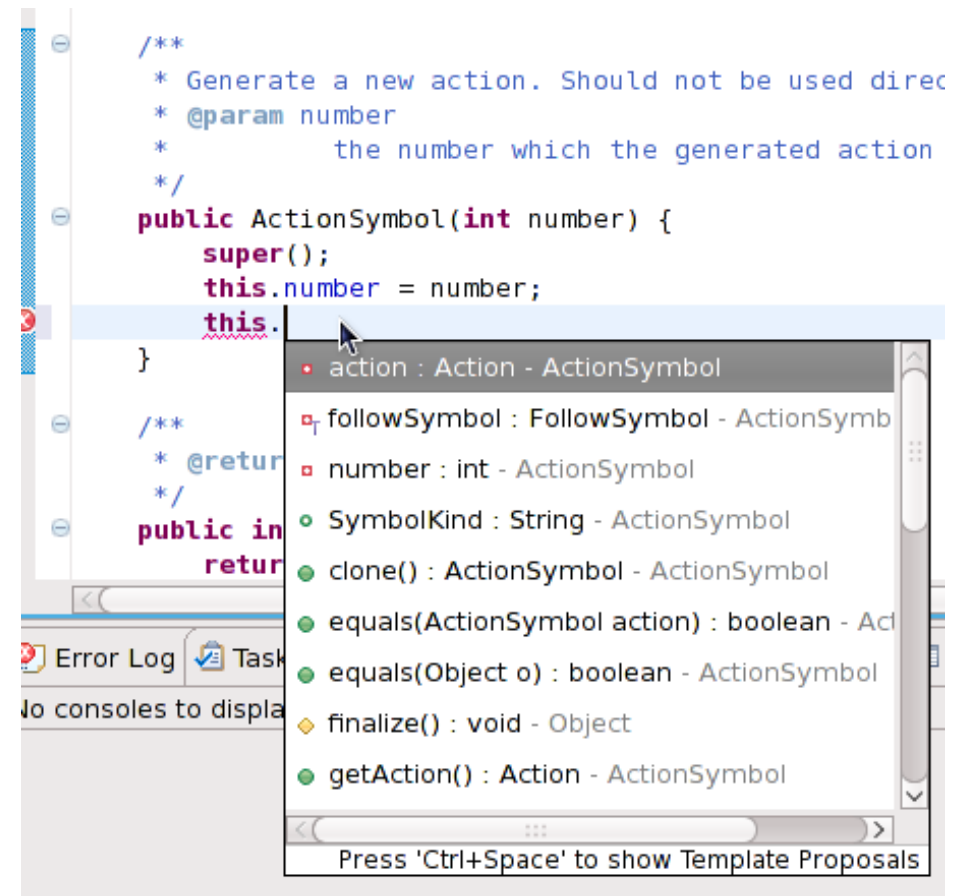
- Incremental Parser
 - Input: An AST and a diff describing the changed input tokens
 - Output: An AST which reuses as much Nodes as possible from the input AST
- Works well for most real life grammars, but yields no general advantage over batch parsers!

Incremental Parsing - Why?

- Why not use it for make?
 - Loading the last AST from disk will probably take longer than parsing from scratch.
 - Creating the diff is non-trivial
- So, why should one write an incremental parser?

Incremental Parsing - Why?

- Programmers need help!
 - This needs an AST to walk on.
 - User expects no delay.
 - Parsing needs to be in „realtime“



```
/**
 * Generate a new action. Should not be used direc
 * @param number
 *         the number which the generated action
 */
public ActionSymbol(int number) {
    super();
    this.number = number;
    this.
}

/**
 * @return
 */
public in
return
```

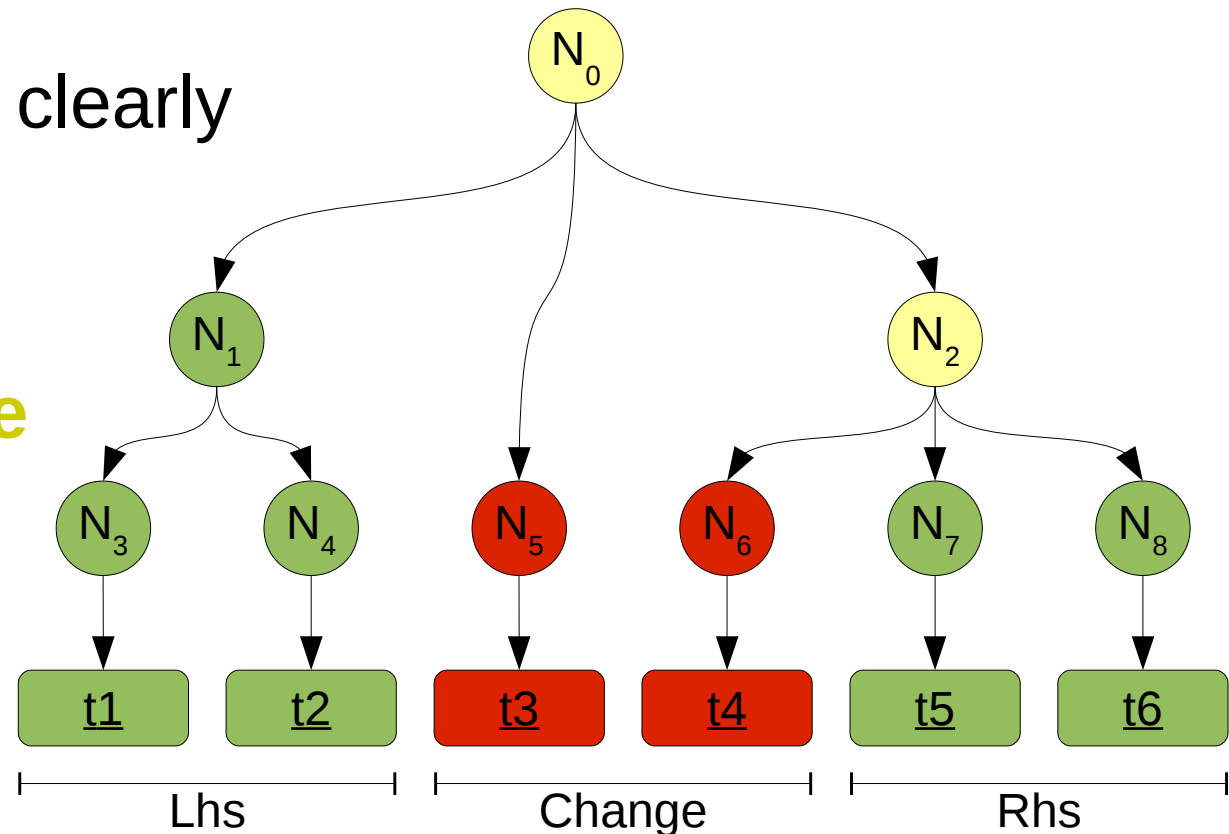
action : Action - ActionSymbol
followSymbol : FollowSymbol - ActionSymbol
number : int - ActionSymbol
SymbolKind : String - ActionSymbol
clone() : ActionSymbol - ActionSymbol
equals(ActionSymbol action) : boolean - ActionSymbol
equals(Object o) : boolean - ActionSymbol
finalize() : void - Object
getAction() : Action - ActionSymbol

Press 'Ctrl+Space' to show Template Proposals

Incremental Parsing - How?

- Tree input:

- Some nodes clearly **reusable**
- Some **not**
- Some **maybe**

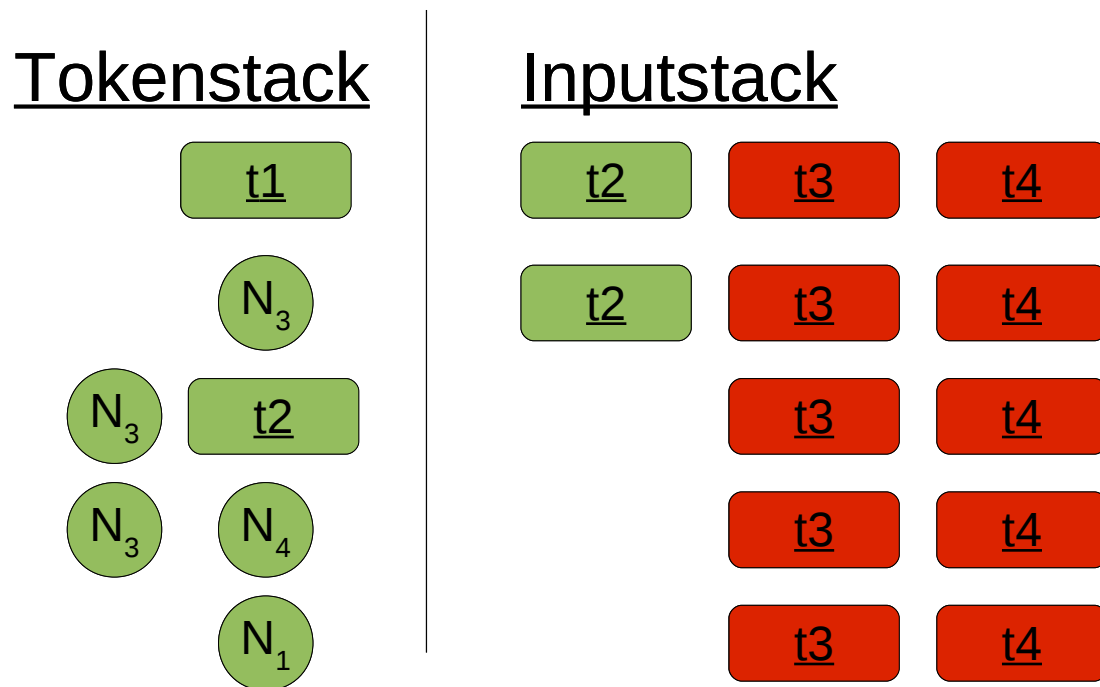


Incremental Parsing - How?

- We will not simply save Parser configurations (although that would be trivial)
- Basic idea: „Short circuiting the parser“
 - We know what the parser did for a given input
- Two Phase design.
 - First phase handle left hand side AST
 - Second phase: handle change spot and right hand side

Incremental Parsing - How?

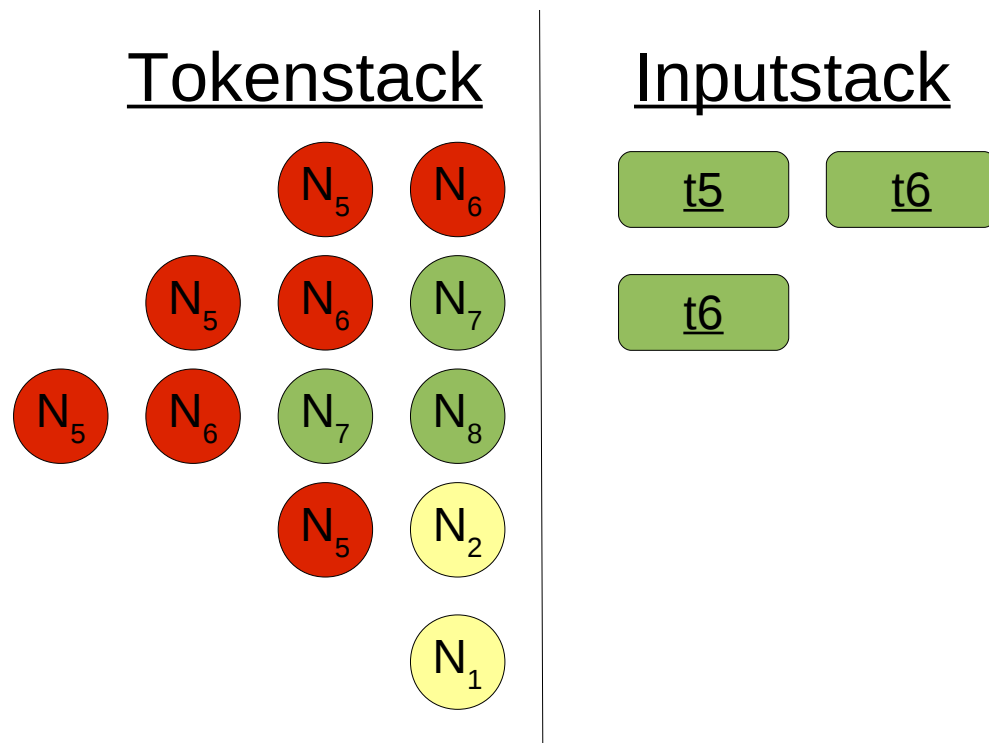
- Batch parsing:



-
- We can directly calculate the final configuration by looking at N_1 , because:
 - It has no changed children
 - The parser is deterministic
 - It is a valid reduction in the current state
 - Phase 1 is:
 - Descend the tree
 - If the Node is unchanged, simulate its reduction
 - else test check its children

-
- Phase 2 is a little bit trickier:
 - Use parent references (new terminals have none)
 - Check for the highest parent node that can be reused
 - (The parse table allows us to check this)
 - This allows reuse of nodes even in new subtrees (!)
 - Use the batch parse steps when needed

- Phase 2:



-
- We assume a balanced AST with n nodes
 - Runtime $O(|\text{change}| + \log^2(n))$:
 - Phase 1 descends the tree: $O(\log(n))$
 - The change is parsed in linear time $O(|\text{change}|)$
 - Phase 2 ascends at least one step after checking $O(\log(n))$ parent nodes: $O(\log^2(n))$
 - Space: None (except parent links, but you'll need them anyway)

- $O(\log^2(n))$ can be dropped to $O(\log(n))$:
 - After reusage directly check for sibling reusage
 - If no sibling, check parents sibling
 - This eliminates unnecessary ascensions
 - Needs some kind of sibling references
 - Sounds cool, but in fact can lead to performance regressions (worst case tree is very unlikely)

-
- Remember the yellow nodes?
 - Why recreate them?
 - We construct bottom up: If one child node already has a parent, reuse it!
 - This would take $O(k)$ time, where k is the maximum amount of child nodes.
 - In PaGe we use Lists, k is not constant
 - But: Only the first and last child need to be checked, since there is only one change

-
- Incremental Parsers are valuable in IDEs
 - For changes with constant length, one can reconcile the AST in $O(\log(n))$ and with no additional space cost
 - This is practical realtime for normal source files (< 64kb)