

A parser for Hack assembly

COMSM1302 Overview of Computer Architecture

John Lapinskas, University of Bristol

Tokens for Hack assembly (reminder)

- **Keywords:**
 - 'A', 'D', 'M',
 - 'JGT', 'JEQ', 'JLT', 'JGE', 'JNE', 'JLE', 'JMP',
 - 'SCREEN', 'KBD', 'SP', 'LCL', 'ARG', 'THIS', 'THAT',
 - and 'R0' through 'R15'.
- **Symbols:** '@', '+', '-', '&', '|', '=', ';', and '!'.
- **Integer literals:** Any base-10 integer in the range 0...32767.
- **Identifiers:** Any string containing no whitespace that's not a keyword and starts with a letter.
- **Newlines.**

EBNF for Hack assembly

There are many equivalent ways to formalise Hack in EBNF. Here's one.

```
⟨instruction⟩ ::= (⟨aInstruction⟩ | ⟨cInstruction⟩), newline;  
⟨aInstruction⟩ ::= “@”, (integerLiteral | identifier | ⟨memoryKeyword⟩);  
⟨memoryKeyword⟩ ::= “SCREEN” | “KBD” | “SP” | “LCL” | “ARG” | “THIS” | “THAT”;
```

EBNF for Hack assembly

There are many equivalent ways to formalise Hack in EBNF. Here's one.

```
⟨instruction⟩ ::= (⟨aInstruction⟩ | ⟨cInstruction⟩), newline;  
⟨aInstruction⟩ ::= “@”, (integerLiteral | identifier | ⟨memoryKeyword⟩);  
⟨memoryKeyword⟩ ::= “SCREEN” | “KBD” | “SP” | “LCL” | “ARG” | “THIS” | “THAT”;  
⟨cInstruction⟩ ::= [⟨assignment⟩], ⟨computation⟩, [⟨jump⟩];
```

EBNF for Hack assembly

There are many equivalent ways to formalise Hack in EBNF. Here's one.

```
⟨instruction⟩ ::= (⟨aInstruction⟩ | ⟨cInstruction⟩), newline;
⟨aInstruction⟩ ::= “@”, (integerLiteral | identifier | ⟨memoryKeyword⟩);
⟨memoryKeyword⟩ ::= “SCREEN” | “KBD” | “SP” | “LCL” | “ARG” | “THIS” | “THAT”;
⟨cInstruction⟩ ::= [⟨assignment⟩], ⟨computation⟩, [⟨jump⟩];
⟨assignment⟩ ::= (“A” | “D” | “M” | (“A”, “D”) | (“A”, “M”) | (“D”, “M”) |
                  (“A”, “D”, “M”)), “=”;
```

EBNF for Hack assembly

There are many equivalent ways to formalise Hack in EBNF. Here's one.

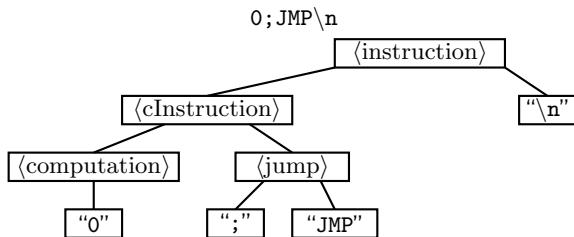
```
⟨instruction⟩ ::= (⟨aInstruction⟩ | ⟨cInstruction⟩), newline;
⟨aInstruction⟩ ::= “@”, (integerLiteral | identifier | ⟨memoryKeyword⟩);
⟨memoryKeyword⟩ ::= “SCREEN” | “KBD” | “SP” | “LCL” | “ARG” | “THIS” | “THAT”;
⟨cInstruction⟩ ::= [⟨assignment⟩], ⟨computation⟩, [⟨jump⟩];
⟨assignment⟩ ::= (“A” | “D” | “M” | (“A”, “D”) | (“A”, “M”) | (“D”, “M”) |
                  (“A”, “D”, “M”)), “=”;
⟨jump⟩ ::= “,”, (“JMP” | “JGT” | “JEQ” | “JLT” | “JGE” | “JNE” | “JLE”);
```

EBNF for Hack assembly

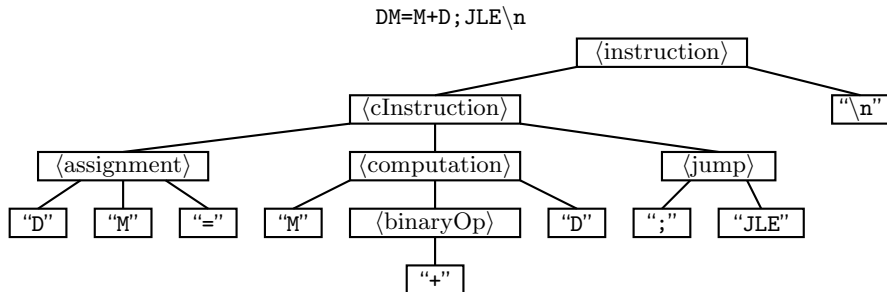
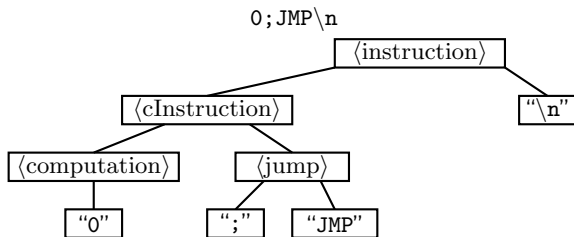
There are many equivalent ways to formalise Hack in EBNF. Here's one.

```
⟨instruction⟩ ::= (⟨aInstruction⟩ | ⟨cInstruction⟩), newline;
⟨aInstruction⟩ ::= “@”, (integerLiteral | identifier | ⟨memoryKeyword⟩);
⟨memoryKeyword⟩ ::= “SCREEN” | “KBD” | “SP” | “LCL” | “ARG” | “THIS” | “THAT”;
⟨cInstruction⟩ ::= [⟨assignment⟩], ⟨computation⟩, [⟨jump⟩];
⟨assignment⟩ ::= (“A” | “D” | “M” | (“A”, “D”) | (“A”, “M”) | (“D”, “M”) |
                  (“A”, “D”, “M”)), “=”;
⟨jump⟩ ::= “;”, (“JMP” | “JGT” | “JEQ” | “JLT” | “JGE” | “JNE” | “JLE”);
⟨computation⟩ ::= “0” |
                ([“-”], “1”) |
                ([“-” | “!”], (“A” | “D” | “M”)) |
                ((“A” | “D” | “M”), (“+” | “-”), “1”) |
                ((“A” | “M”), ⟨binaryOp⟩, “D”) |
                (“D”, ⟨binaryOp⟩, (“A” | “M”));
⟨binaryOp⟩ ::= “+” | “-” | “&” | “|”;
```

Example CSTs for Hack assembly



Example CSTs for Hack assembly



How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

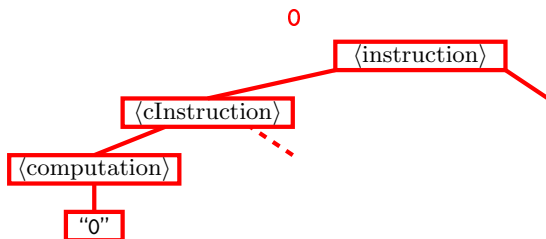
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

0

How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

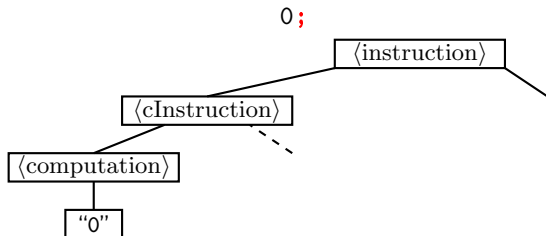


0 can arise only as an entire $\langle \text{computation} \rangle$, which must be in a $\langle \text{cInstruction} \rangle$, which must be in an $\langle \text{instruction} \rangle$.

We also know that this $\langle \text{cInstruction} \rangle$ has no $\langle \text{assignment} \rangle$, although it might have a $\langle \text{jump} \rangle$.

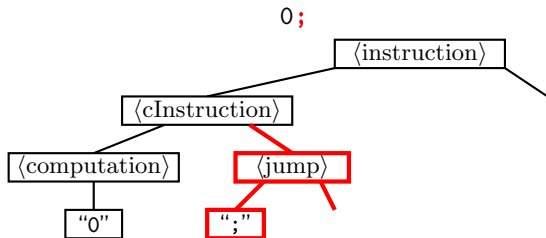
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

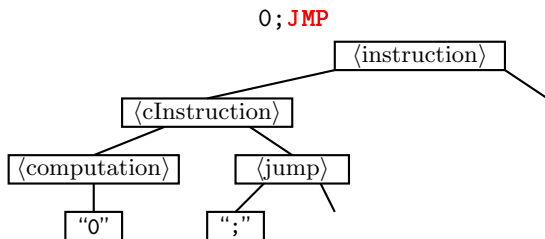
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



; can arise only as the first term of a $\langle \text{jump} \rangle$, which must be part of the $\langle \text{cInstruction} \rangle$ we already added.

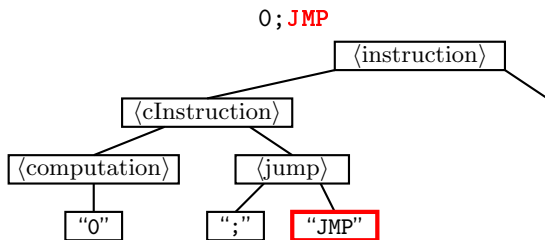
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

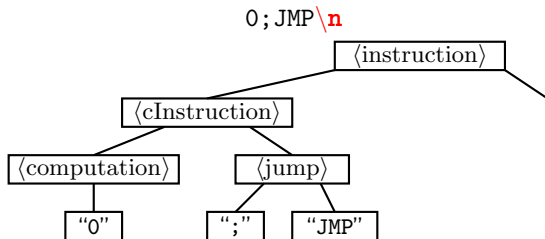
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



JMP must be the continuation of that <jump>.

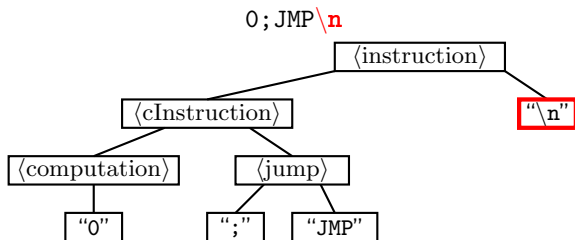
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

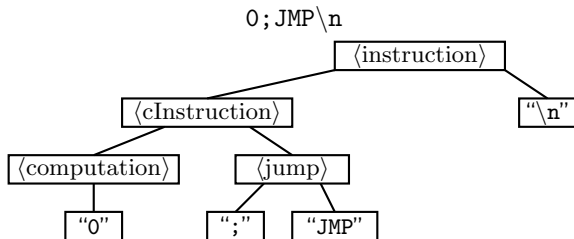
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



Finally, \n must be the end of the <instruction>.

How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



In each case, if there had been no way to fit the next token into our existing CST (e.g. on parsing the 1 in 01; JMP), we would return an error.

How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

D

How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

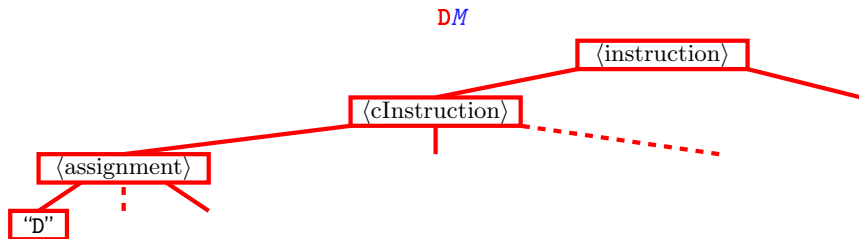
DM

We can't tell immediately how *D* should fit into the CST. It could be part of either an $\langle \text{assignment} \rangle$ or a $\langle \text{computation} \rangle$.

We can check which by looking ahead to the *second* token. If it's *M* or *=*, then we must be in an $\langle \text{assignment} \rangle$. If it's *+*, *-*, *&*, *|*, *;* or $\backslash n$ we must be in a $\langle \text{computation} \rangle$.

How do we build a Hack CST?

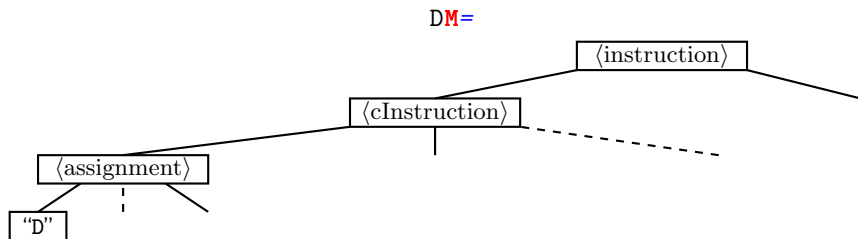
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



In this case, the next token is M. So D must be the start of an $\langle \text{assignment} \rangle$, which might also contain an M and will contain an =.

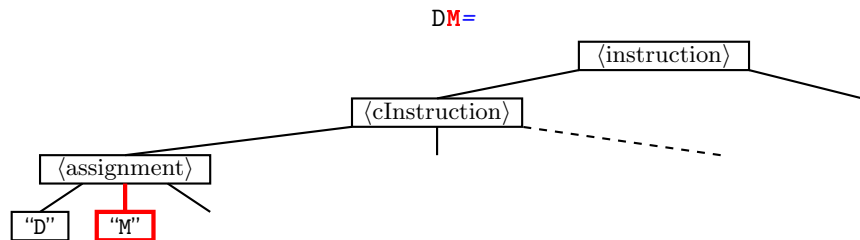
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

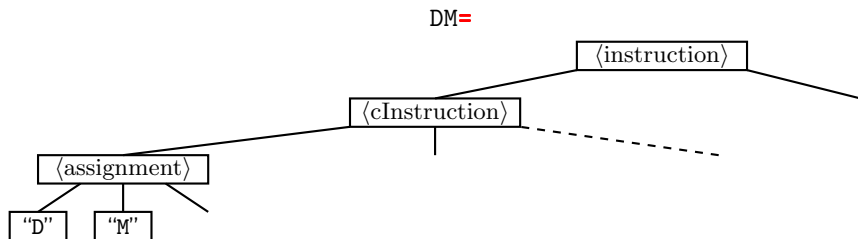


Again, looking ahead from M we see =, so we know it must be in an $\langle \text{assignment} \rangle$ as expected. (We could also use the fact that we know the next term must be part of an $\langle \text{assignment} \rangle$ to fit into the existing CST.)

The $\langle \text{assignment} \rangle$ must be part of a $\langle \text{cInstruction} \rangle$, which must be in an $\langle \text{instruction} \rangle$.

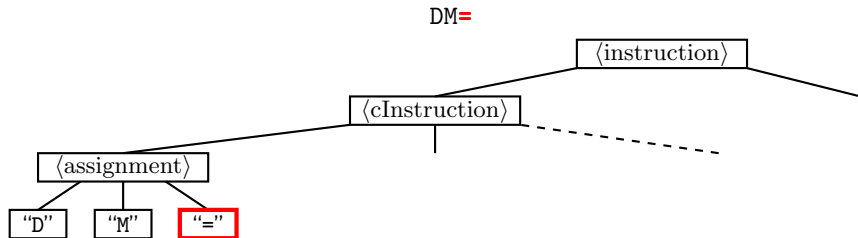
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

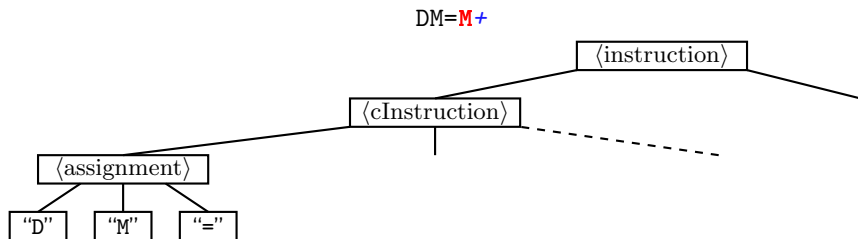
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



= can only fit into that `<assignment>`.

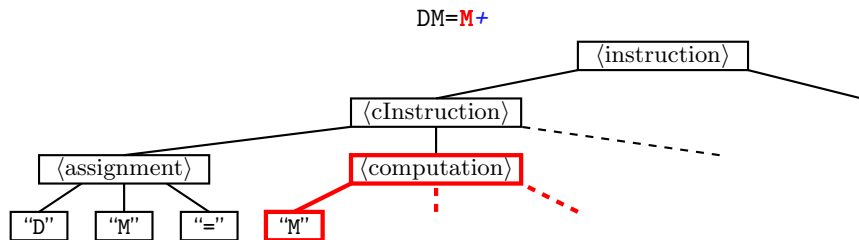
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

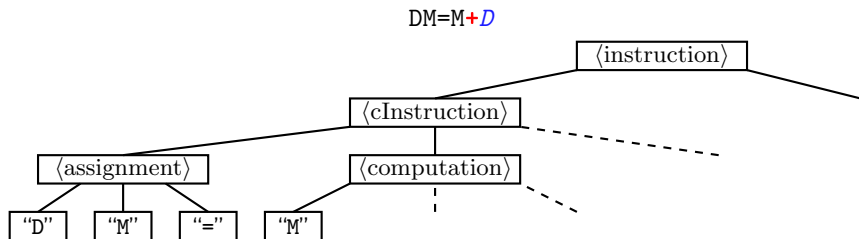


For this M, we again look ahead to see that the next character is +, so we must be in a $\langle \text{computation} \rangle$ as expected.

The $\langle \text{computation} \rangle$ may or may not contain a $\langle \text{binaryOp} \rangle$ and a second term.

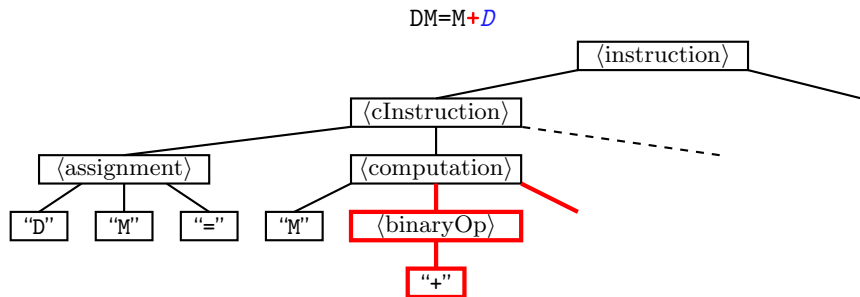
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.

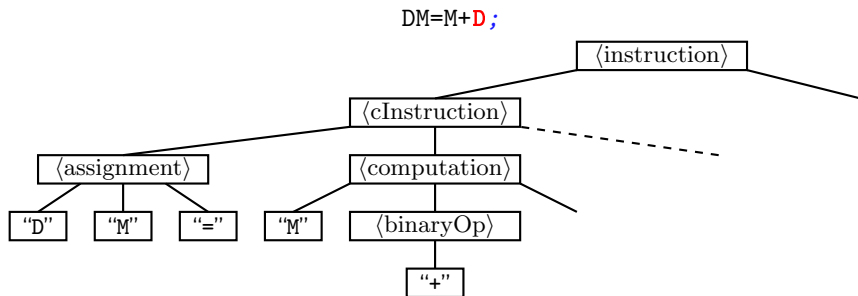


+ must be part of a $\langle \text{computation} \rangle$, but might be of the $M+1$ form or of the $M+D$ form (which leads to a different CST since the + is a $\langle \text{binaryOp} \rangle$).

Again, by looking ahead another token, we can tell that it is a $\langle \text{binaryOp} \rangle$. This also tells us we should expect the next term of the $\langle \text{computation} \rangle$ to be a D.

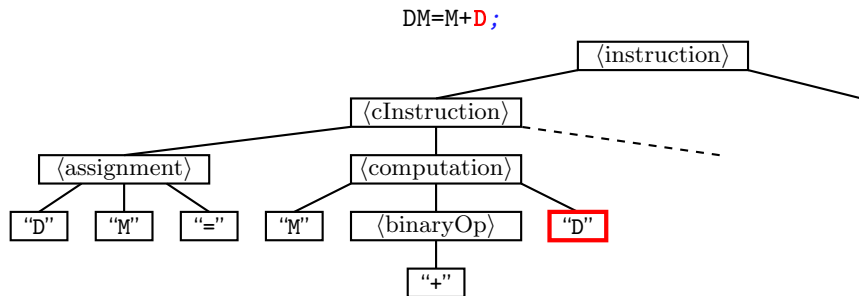
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

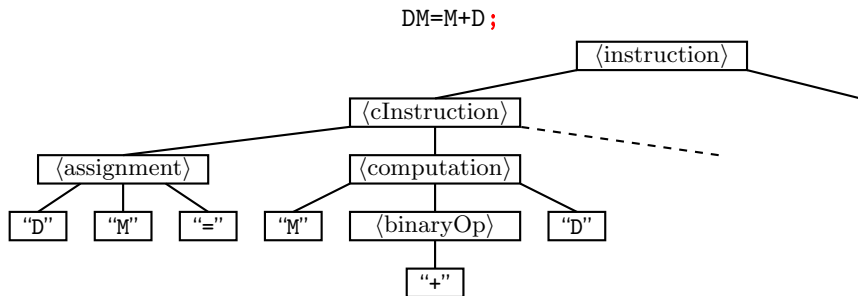
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



As expected, we see a D. Looking ahead we see a ;, so this is part of a \langle computation \rangle as expected.

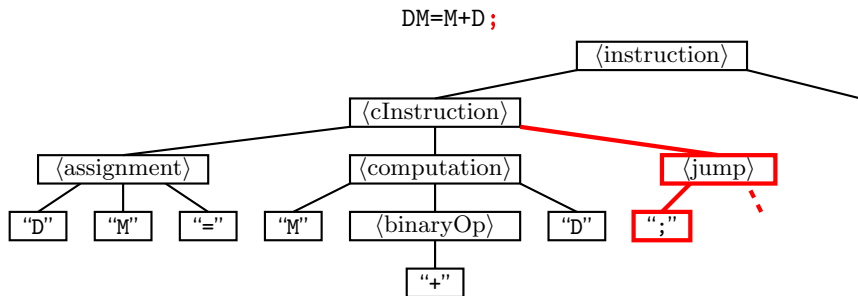
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

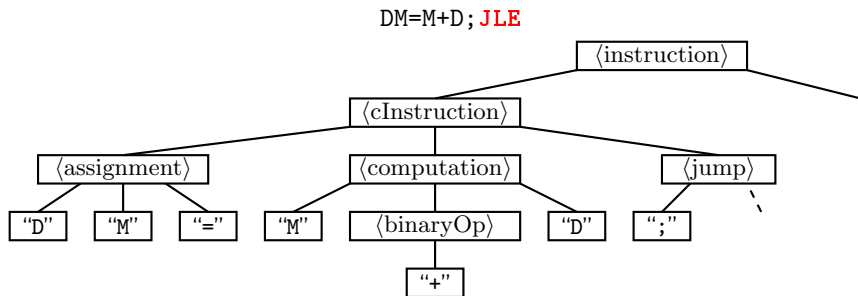
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



As in the previous example, ; must be part of a <jump>...

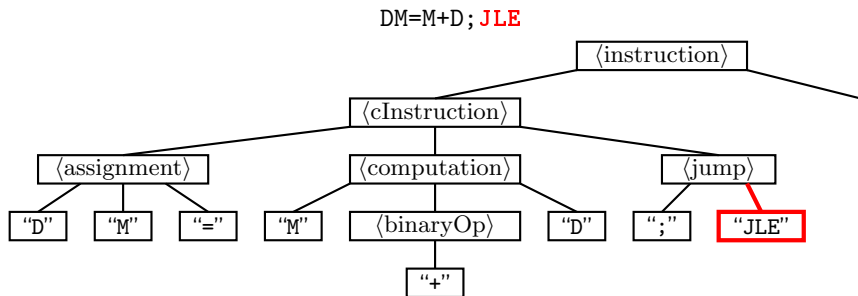
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

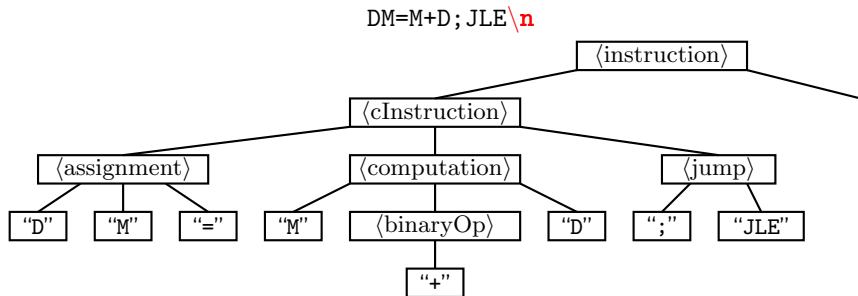
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



which is completed by this JLE...

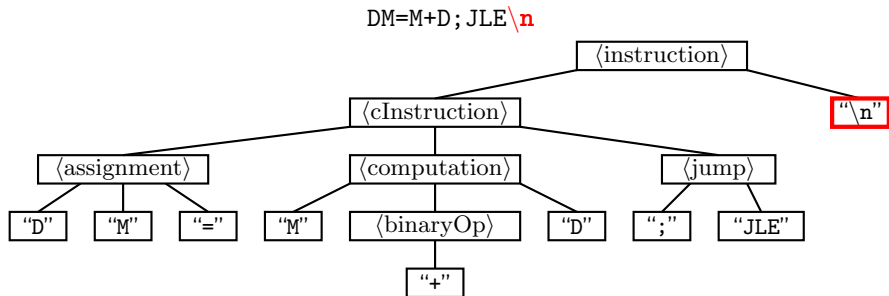
How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



How do we build a Hack CST?

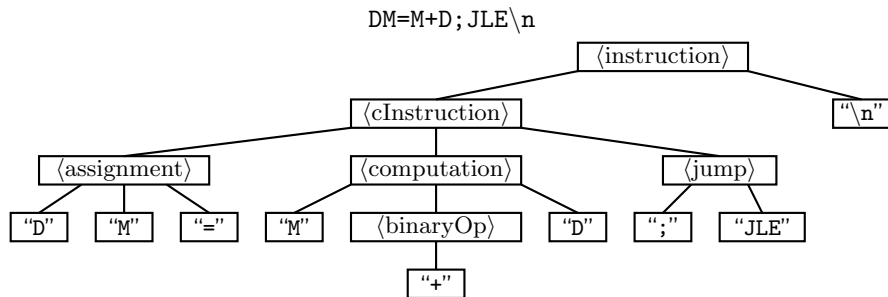
In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



leaving only the final newline.

How do we build a Hack CST?

In this unit, we will only consider **LL parsing**. We go through tokens from left to right, building the CST from the top down by looking at only the next few tokens.



Grammars which can be LL-parsed looking by looking k tokens ahead are called **LL(k)**. For example, Hack is LL(2) (because of e.g. DM=M+D) but not LL(1).

Warning: I'm sweeping some very deep rabbit-holes under the rug here. Not all grammars are LL(k), and this isn't even the formal definition of LL(k). But for a quick and dirty parser implementation it's good enough!

How do we actually parse Hack?

CSTs are extremely useful for graceful error-handling and for compiling complex expressions like `i=(k+j++)/m`. But here, they're overkill.

What we actually need to know:

- Is the statement an A-instruction or a C-instruction?
- If it's an A-instruction, what value should we load into A?
- If it's a C-instruction, what are the values of `dest`, `comp`, and `jump`?

How do we actually parse Hack?

CSTs are extremely useful for graceful error-handling and for compiling complex expressions like $i=(k+j++)/m$. But here, they're overkill.

What we actually need to know:

- Is the statement an A-instruction or a C-instruction?
 - Is the first token an @?
- If it's an A-instruction, what value should we load into A?
- If it's a C-instruction, what are the values of dest, comp, and jump?

How do we actually parse Hack?

CSTs are extremely useful for graceful error-handling and for compiling complex expressions like $i=(k+j++)/m$. But here, they're overkill.

What we actually need to know:

- Is the statement an A-instruction or a C-instruction?
 - Is the first token an @?
- If it's an A-instruction, what value should we load into A?
 - What comes after the @?
 - If it's an identifier, what RAM/ROM address is it? (Via symbol tables.)
- If it's a C-instruction, what are the values of `dest`, `comp`, and `jump`?

How do we actually parse Hack?

CSTs are extremely useful for graceful error-handling and for compiling complex expressions like $i=(k+j++)/m$. But here, they're overkill.

What we actually need to know:

- Is the statement an A-instruction or a C-instruction?
 - Is the first token an @?
- If it's an A-instruction, what value should we load into A?
 - What comes after the @?
 - If it's an identifier, what RAM/ROM address is it? (Via symbol tables.)
- If it's a C-instruction, what are the values of `dest`, `comp`, and `jump`?
 - Does = appear, and which of A, D and M appear to the left of it?
 - Does ; appear, and which jump instruction appears to the right of it?
 - What's to the right of the = and the left of the ;? (30-ish possibilities.)

How do we actually parse Hack?

CSTs are extremely useful for graceful error-handling and for compiling complex expressions like $i=(k+j++)/m$. But here, they're overkill.

What we actually need to know:

- Is the statement an A-instruction or a C-instruction?
 - Is the first token an @?
- If it's an A-instruction, what value should we load into A?
 - What comes after the @?
 - If it's an identifier, what RAM/ROM address is it? (Via symbol tables.)
- If it's a C-instruction, what are the values of `dest`, `comp`, and `jump`?
 - Does = appear, and which of A, D and M appear to the left of it?
 - Does ; appear, and which jump instruction appears to the right of it?
 - What's to the right of the = and the left of the ;? (30-ish possibilities.)

We can answer all these questions with simple logic on tokens.

(Making the CST is still a good exercise if you have extra time, though!)

The Hack assembler: A summary

Pass 1: Lexing. For each line:

- Remove any comments and whitespace.
- If the line is empty, skip it.
- If the line is a label, add it to the symbol table along with the ROM address corresponding to the current line.
- Otherwise, break the line into tokens and output to a temporary file.

Pass 2: Parsing. For each \langle instruction \rangle (separated by newline tokens):

- If the \langle instruction \rangle starts with an '@' token:
 - If it uses a new variable, allocate RAM and add it to the symbol table.
 - If it uses an existing variable or label, retrieve the RAM/ROM address for it from the symbol table.
 - Generate and output the corresponding A-instruction.
- Otherwise:
 - Break it down into an assignment, a computation, and a condition.
 - Map these to appropriate values of `dest`, `comp` and `jump` respectively.
 - Generate and output the corresponding C-instruction.