

14. Context-Sensitive Languages

Goal: • Introduce context-sensitive languages (CSLs) and Turing machines.

- Prove decidability of membership in CSLs.
- Show correspondence of CSLs and languages accepted by linear-bounded Turing machines.
- This also establishes equivalence of the recursively enumerable languages and the languages accepted by Turing machines.

Definition:

• A type-0 grammar is a grammar $G = (N, \Sigma, P, S)$ with $P \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$.

• The grammar is context-sensitive (or type-1), if for all productions

$$\alpha_1 \rightarrow \alpha_2 \in P \quad \text{we have } |\alpha_2| \leq |\alpha_1|.$$

Moreover, we may have $S \rightarrow \epsilon$, but then there is no rule that has S in its right-hand side.

• A language $L \subseteq \Sigma^*$ is context-sensitive,

if there is a context-sensitive grammar G with $L = L(G)$.

• A language $L \subseteq \Sigma^*$ is recursively enumerable,

if $L = L(G)$ for G of type-0.

Example:

$$L = \{w \in \{a,b,c\}^* \mid \#_a(w) = \#_b(w) = \#_c(w)\}.$$

We have $L = L(G)$ with

$$G = S \rightarrow \epsilon \mid R$$

$$R \rightarrow ABC \mid ABCR$$

$$A \rightarrow a \quad B \rightarrow b \quad C \rightarrow c$$

$$AB \rightarrow BA$$

$$AC \rightarrow CA$$

$$BA \rightarrow AB$$

$$BC \rightarrow CB$$

$$CA \rightarrow AC$$

$$CB \rightarrow BC.$$

Let $L(G)$ be a context-sensitive language (over Σ).

The membership problem is

MEMBERSHIP $L(G)$:

Given: $w \in \Sigma^*$.

Question: $w \in L(G)$?

Theorem: For every context-sensitive grammar G ,
MEMBERSHIP $L(G)$ can be solved in $2^{O(|w|)}$.

Proof: Since the productions are length preserving,
once we derived a sentential form of length $> |w|$,
there is no chance to derive w from it.

• We thus enumerate all sentential forms of length $\leq |w|$
and see whether we find w .

• There are $(|N| + |\Sigma| + 1)^{|w|} = 2^{O(|w|)}$ many,

which yields the time bound.

Note that N and Σ are not part of the input

and therefore $|N| + |\Sigma| + 1 = 2^c$

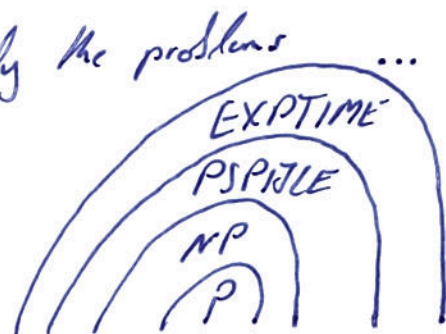
for c a constant that is independent of the input (word w). \square

The problem is hard:

• The correspondence with linear-bounded Turing machines
that we establish below shows that

$$CSL = PSPACE,$$

the context-sensitive languages are precisely the problems
solvable with polynomial space.



- Since there are PSPACE-hard problems, there are context-sensitive languages $L(G)$ for which membership is PSPACE-hard.

→ PSPACE is an important complexity class for verification problems:

- ↳ Reachability in Boolean Nite-programs
- ↳ Reachability in multi-threaded programs.

→ Context-sensitive aspects show up in compilers

- ↳ Parameters in procedures
- ↳ Scopes of objects.

14.1 Turing Machines

Goal: Define an automaton model, the Turing machine (TM), for context-sensitive languages and for recursively-enumerable languages.

Idea: Drop the last-in-first-out restriction for pushdowns and access arbitrary information stored about the computation so far.

- For context sensitivity, we just bound the amount of information we can store.

Church's thesis: Turing's idea was more than capturing a class of languages. He wanted to formulate the idea of computability itself.

(Alonzo Church wanted to capture what is an algorithm.)

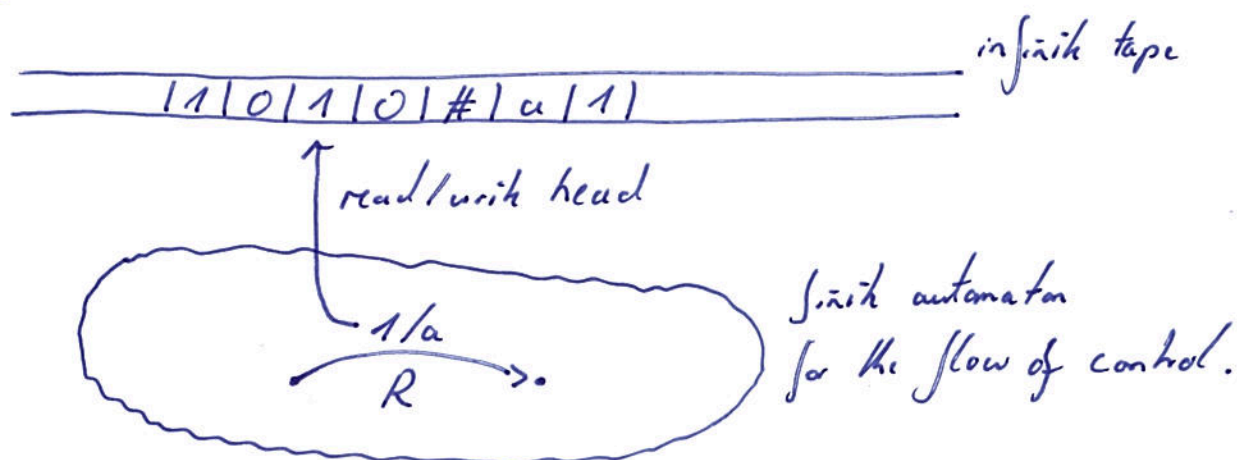
(Church 1936) Up to today, we believe he succeeded.

EVERY problem for which we found an algorithm (C, C++, Java) we were able to solve with a Turing machine (even with low overhead).

Church's thesis states that it will stay this way:

Everything that is computable is computable by a TM.

Technically:



Definition (Syntax of a TM):

A Turing machine (TM) is a tuple $M = (Q, \Sigma, T, q_0, \sqcup, \delta, Q_f)$

with

- Q a finite set of states, $q_0 \in Q$ the initial state,
- $Q_f \subseteq Q$ the set of final states,
- Σ the input alphabet,
- T the tape alphabet with $\Sigma \subseteq T$ and $\sqcup \in T \setminus \Sigma$ the blank symbol,
- $\delta : Q \times T \rightarrow T \times \{L(\text{left}), R(\text{right}), N(\text{central})\} \times Q$,
in which case the TM is deterministic (DTM), or
 $\delta \subseteq Q \times T \times T \times \{L, R, N\} \times Q$ for a non-deterministic TM (NTM)

To define the semantics of a TM,
we need the notion of a configuration, the state of the TM at runtime.

Definition (Semantics of a TM):

Let $M = (Q, \Sigma, T, q_0, \sqcup, \delta, Q_f)$ be a TM.

- A configuration of M is a word from $T^* Q T^*$.

If $uqav$, then uav are the cells of the tape already visited,
 q is the current state and

$a \in \Sigma$ is the symbol currently under the read-write head.

Given input $x \in \Sigma^+$, the corresponding initial configuration is $q_0 x$.

Given input ϵ , the initial configuration is $q_0 \sqcup$.

• A configuration is accepting if it is from $T^* Q_F T^+$.

• The transition relation among configurations

$$\rightarrow \subseteq (T^* Q T^+)_x (T^* Q T^+)$$

is defined by

$$u a q b v \rightarrow u q' a c v, \quad \text{if } q \xrightarrow[L]{b/c} q' \in \delta$$

$$u a q b v \rightarrow u a q' c v, \quad \text{if } q \xrightarrow[N]{b/c} q' \in \delta$$

$$u a q b v \rightarrow u a c q' v, \quad \text{if } v \neq \epsilon \text{ and } q \xrightarrow[R]{b/c} q' \in \delta$$

$$u q b \rightarrow u c q' \sqcup, \quad \text{if } q \xrightarrow[R]{b/c} q' \in \delta$$

$$q b \sqcup \rightarrow q' \sqcup c \sqcup, \quad \text{if } q \xrightarrow[L]{b/c} q' \in \delta$$

We use \rightarrow^* for the reflexive and transitive closure of \rightarrow .

• The language of M is

$$L(M) := \{ w \in \Sigma^* \mid q_0 w \rightarrow^* u q v \in T^* Q_F T^+ \} \\
(q_0 \sqcup \text{ for } \epsilon)$$

For the characterization of the context-sensitive languages,
we define Turing machines that only use the part of the tape
given by the input (plus a left and a right end marker).

Definition (Linear-bounded automaton):

• A linear-bounded automaton (LBA) is an NTM

$$M = (Q, \Sigma, \Sigma \cup \{ \$, \$ \}, q_0, \delta, Q_F).$$

Instead of a blank \sqcup , we have a left end marker $\$l$
and a right end marker $\$r$.

Moreover, there are two restrictions

- 1) There are no moves left from $\$l$ and right from $\$r$.
- 2) The symbols $\$l$ and $\$r$ are never overwritten.

The language of M is

$$L(M) := \{ w \in \Sigma^* \mid q_0 \$l w \$r \xrightarrow{*} u q_f v \in T^* Q_f T^* \}$$

A tape compression result in complexity theory shows that an amount of tape linear in the size of the input leads to the same computational power (as restricting to the input itself). Hence the name linear-bounded automaton.

14.2 Equivalence with Context-Sensitive Languages

Goal: Show that the languages accepted by LBAs are precisely the context-sensitive languages.

Theorem (Kuroda '64):

A language $L \subseteq \Sigma^*$ is accepted by an LBA iff it is context-sensitive.

Proof:

\Leftarrow Let $L = L(G)$ with $G = (N, \Sigma, P, S)$.

We give an NTM M that accepts L .

The idea is to apply the productions backwards until the start symbol is found.

\hookrightarrow Consider an input $x \in \Sigma^*$.

$\hookrightarrow M$ non-deterministically selects a production

$$\alpha \rightarrow \beta \in P$$

\Leftarrow and finds an arbitrary occurrence of β in x .

↳ Now β is replaced by α ,
potentially moving letters to the left of $|\alpha| < |\beta|$.

↳ If we reach S , we stop and accept.

Otherwise, we repeat the above non-deterministic procedure.

Since the productions are length preserving ($|\beta| \geq |\alpha|$),
we never exceed the part of the tape used by the input.

⇒ Let $L = L(M)$ for an LBA $M = (Q, T, \Sigma \cup \{l_e, l_r\}, q_0, \delta, Q_f)$.

Idea: • We give a grammar G that operates on sentential forms
representing the configurations of M .

• We will have to generate the initially given word
from these sentential forms.

Therefore, they cannot be longer than the input (we cannot delete).

• Note that this bound works for the intermediary configurations
due to the linear boundedness assumption.

Technically:

• We will store pairs

$(a_1, a_1) \dots (a_n, a_n)$

where $a_1 \dots a_n$ forms the current tape content and
 $a_1 \dots a_n$ is the initial input.

• The letters a_i are taken from the alphabet Δ with
 $\Delta' := T \cup (\{l_e, l_r\} \times T)$ and $\Delta := \Delta' \cup (Q \times \Delta') \cup (\{l_e, l_r\} \times Q \times T)$.

So an intermediary configuration

$l_e \times q \times a \times l_r$

derived from input

$S_c a^s a^t S_r$

has the shape

$$((S_c, x), a) ((q, y), b) (a, a) ((S_r, z), b).$$

- It is immediate to simulate the moves of the LBA by considering only the current configuration (one has to trick a bit at the border (*)).

So if $q \xrightarrow{a/b} q' \in S$, then we obtain the productions

$((q, a), x) \cdot (c, y) \rightarrow (b, x) \cdot ((q', c), y)$ for all $x, y \in \Sigma$ in the grammar.

(*) For the border, we use

(q, S_c, a) to indicate that the head is on the left end marker

and (S_c, q, a) to indicate that the head is on the first letter.

Construction:

We define

$$G := (N, \Sigma, P, S)$$

with $N := \{S, A\} \cup (\Delta \times \Sigma)$

$$P := \{ S \rightarrow A \cdot ((S_r, a), a) \mid a \in \Sigma \}$$

$$\cup \{ A \rightarrow A \cdot (a, a) \mid a \in \Sigma \}$$

$$\cup \{ A \rightarrow ((q_0, S_c, a), a) \mid a \in \Sigma \}$$

\cup Productions that simulate the behavior of M

$$\cup \{ ((q_s, a), b) \rightarrow b \mid q_s \in Q_s, a \in \Delta', b \in \Sigma \}$$

$$\cup \{ (a, b) \rightarrow b \mid a \in \Delta', b \in \Sigma \}$$

Note: The construction works the same way

if we drop

- ↳ the length-preserving condition for grammars and
- ↳ the linear boundedness assumption for TMs.

Corollary:

The languages accepted by an NTM are precisely the recursively enumerable languages.

14.3 Determinism

Recall: For finite automata, we know that NFAs and DFAs accept the same languages.

For pushdown automata, we know that DPDAs accept a strict subclass of the context-free languages.

Open: Kuroda, in his 1964 paper, posed two questions:

(1) The first question is indeed on

the role of determinism for LPAs:

Are the languages accepted by NLPAs

precisely the languages accepted by DLPAs?

It is still open!

(2) The second question was on complementation:

Are the languages accepted by NLPAs

(the context-sensitive languages)

closed under complement?

Kuroda showed that $\neg(2) \Rightarrow \neg(1)$.

This did not help much:

Immerman & Szepietowski proved (2) to hold!

We will see the proof in the next lecture.

Goal: Show that for TM, NTMs and DTMs accept the same languages.

Theorem: L is accepted by an NTM iff L is accepted by a DTM.

Proof:

" \Leftarrow " \forall Every DTM is an NTM.

" \Rightarrow " The idea is to traverse the computation tree of the NTM in breadth-first manner.

• Let $L = L(M_1)$ with M_1 an NTM.

For any control state and any tape symbol, there is a finite number of choices for the next move of M_1 .

Let $r \in \mathbb{N}$ be the maximal number of choices for any pair of control state and tape symbol.

Then any finite sequence of choices can be represented by a sequence of digits from 1 to r .

Not all such sequences may represent choices of moves, there may be fewer than r choices in some situations.

• We simulate M_1 by a DTM M_2 that has three tapes.

The three tapes can be compressed into one tape

by taking letters to be triples $\begin{pmatrix} a \\ b \\ c \end{pmatrix}$ or

by writing the tapes one after the other.

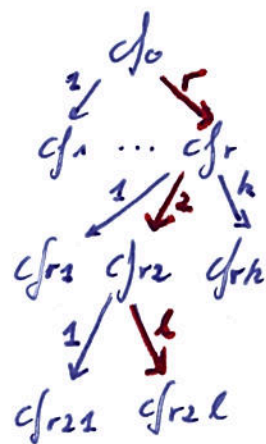
The first tape of M_2 holds the input.

On the second tape, M_2 generates sequences of the digits 1 to r .

in a systematic manner that implements a breadth-first search:

• The sequences are generated in increasing length, starting from the shortest.

• Sequences of equal length are generated in numerical order.



The path is given by the sequence $r, 2, l$

// We will discuss such techniques when we come to computability.

- For each sequence generated on tape 2, M_2 copies the input onto tape 3 and then simulates M_2 on tape 3 — using the sequence on tape 2 to resolve non-deterministic choices.
If there is a sequence of choices leading to acceptance of M_2 , it will eventually be generated on tape 2.
When M_2 then simulates M_2 , it will accept.
If no sequence of choices of M_2 leads to acceptance, M_2 will not accept. □

Note: The construction does not carry over to LBAs as it is not clear that it preserves linear boundedness (the sequence of choices may be long).