

A language class is a set of languages

Each language class has two properties

- Decision properties
- Closure properties

Usually we can use algorithms to find these properties.

Closure properties says that given a language in a class, an operation produces another language of the same class.

A language can be described in formal or informal way (such as $2^{\{0,1\}}$ is a language) (defined by PE or by FA)

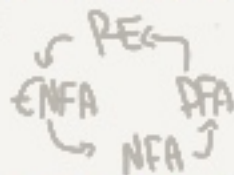
You cannot answer a question about a language unless you have a formal description.

Thus a decision property is an algorithm that takes a formal description and tells whether some property holds

Property such as is it finite? Can it fail? Things that you cannot do with programs.

The membership property takes a string and tells whether it is accepted by the DFA.

The algorithm simulates the DFA on the input. If the regular language is represented in an NFA or RE, we convert to DFA by following the circle of conversions



The emptiness property, given a regular language, does it contain any string at all?

- Assume it is represented by DFA
- Find reachable states by using BFS or DFS on the start state.
- If at least one final state is marked, then it accepts or least one input, else it is empty

The infiniteness problem

- Start with DFA of the language
- **Key Idea** If the DFA has n states, and the language contains only strings of length n or more then the language is infinite
- Otherwise is finite to n or less

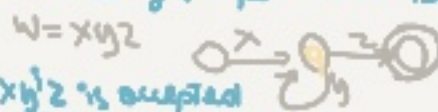
10.06 on...

Proof of **key idea** "If the DFA accepts any string whose length is at least the number of states n , then it accepts an infinite number of strings"

- First observe that a string of length n or more has at least $n+1$ states along its path



If there are only n states, and there are $n+1$ states along the path, then two states along the path must be the same



Then xy^iz is accepted by the DFA for any $i \geq 0$

Since y is not ϵ , the DFA can accept an infinite number of strings.

We still don't have an algorithm because we can't test the infinite number of strings of length $> n$

OTHER KEY IDEA It is sufficient to test strings between n and $2n-1$



- y is the first cycle, so xy cannot be greater than n . Some state in the first $n+1$ states solely repeats.
- If w is of length $2n$ or more, there is still a shorter string of length $\leq n$
- Keep shortening to reach $[n, 2n-1]$

The algorithm tests for membership all strings of length between n and $2n-1$. If it is accepted is infinite.

...

A more efficient algorithm:

- Eliminate states not reachable from start and final state
- Check for cycles in the remaining graph.
- If you don't find cycles, the language is finite

PUMPING LEMMA number of states

For every language L , there is an integer n such that for every string w in L of length $\geq n$ we can write $w = xyz$ such that

- $|xy| \leq n$
- $|y| > 0$
- For all $i \geq 0$, xy^iz is in L

y is the label of the first subthing of w that goes from the start to the same state

The prefix of xy in w is at most n , come to the label of the first cycle in the DFA

- y is not ϵ because it connects two different occurrences of the same state along path of w

- xy^iz is in L for $i \geq 0$



Proof for containment Given L and M , is $L \subseteq M$?

- Algorithm uses the product construction.

- We define Q as Q_1 but not Q_2

Given a DFA, find the equivalent DFA with the fewest states

- Since we can test equivalence, given a DFA we can find the DFA with fewest states accepting $L(A)$

Algorithm, creates a table with pairs of states and figure out which pairs are distinguishable (leads one to a final state and the other to a non).

Proof for equivalence

- Given regular languages L and M is $L = M$?

- Algorithm involves constructing the product DFA for L and M

- Suppose have set of states Q and R

- Product has set of states $Q \times R$ pairs (q, r) q from Q , r from R

- Start (q_0, r_0)

- Transitions $\delta([q, r], a) = [\delta_L(q, a), \delta_M(r, a)]$

- Make the final state be all the pairs such that one is a final state and the other isn't. If w reaches only one final state, the two languages are not the same

STATE MINIMIZATION

Basis, Mark pairs with one final state

Induction, mark $[q, r]$ if for some input a , $[q(a), r(a)]$ is marked

After no more marks are possible, the unmarked pairs are equivalent and can be merged into one state

CLOSURE PROPERTIES

Given an operation between two languages, when the argument are of a class, so is the result

- L M are regular, then $L \cup M$ is also regular
Let R and S be the RE defined, then $R \cup S$ is a RE whose language is $L \cup M$

- RS is a RE whose language is LM **CONCAT.**
- R^* is a RE whose language is R^* **KLEENE**

- If L and M are regular, then so is $L \cap M$
"We cannot use REs very easily to prove this one"
- construct C , the product automaton of A and B
- make the final states be the pairs consisting of final states of both A and B

If L and M are regular, then so is $L - M$

- construct C ...
- make the final states be the pairs where A -state is final and B -state is not

The complement of a language L is $\Sigma^* - L$

- Since we know that regular languages are closed and start and difference, we know that complement is also closed.

HOMOMORPHISM

They are transformations on symbols which replace each symbol by a string

When a homomorphism is applied to a regular language, the result is a regular language

i.e. $h(a) = ab, h(b) = a$
Apply it to any string in order and concatenate the results
 $h(01010) = ababab$

Inverse homomorphism

$h^{-1}(L)$
The result is the set of strings w such that $h(w)$ is in L

CONTEXT FREE

Is a class of languages

- Applications are processing natural language, and computer languages
They're essential when designing the parser of a language (in a compiler it puts tokens together into the proper structure)

Context free languages are defined by context-free grammars

Every regular language has a context free grammar describing it, not all context free languages are also regular

Many of these non-regular languages, are languages that involve nested structures

The central elements of the grammar are variables, symbols that generate particular sets of strings

The variables, or the set of strings they generate, are defined recursively, in terms of one another

The rules that defines the language of the variables are called productions
It involves concatenation $A \rightarrow xy$ the concatenation of $L(x)$ and $L(y)$ is a subset of $L(A)$

A language may have several productions, it is defined as the union of the languages described by the right side of each production

eg. CFG for $\{0^n 1^n | n \geq 1\}$

productions $S \rightarrow 01$ BASIS RULE
 $S \rightarrow 0S1$ INDUCTIVE RULE

~ 01 is in the language
 \sim If w is in the language, then so is $0w1$

CFG FORMALISM

Terminology

They are analogous to the input symbols of an automaton. They form the alphabet for the language being defined

Variables are non-terminals. They are something like the states of an automaton, but more powerful. The **start symbol** variable it is the language of this variable that the grammar defines, and the other variables are used as auxiliaries

The productions of the grammar, which are akin to the transition function of an automaton, have the form variable string of (head) \rightarrow variables and terminals (body)

CONVENTIONS

A, B, C, \dots and S are variables

a, b, c, \dots are terminals

\dots, X, Y, Z are either terminals or variables

\dots, w, x, y, z are strings of terminals only

$\alpha, \beta, \gamma, \dots$ are strings of terminals and/or variables

eg. CFG for $\{0^n 1^n | n \geq 1\}$

terminals $0, 1$
variables $\{S\}$ **start symbol**
productions $S \rightarrow 01 \quad S \rightarrow 0S1$

A derivation is a sequence of strings that typically have both terminals and variables.

We say $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if $A \rightarrow \gamma$ is a production

Iterated derivation

\Rightarrow^* means zero or more derivation steps

Means $\alpha \Rightarrow^* \beta$ if for any string α and β such that $\alpha \Rightarrow \beta$ then $\alpha \Rightarrow^* \beta$

SENTENTIAL FORMS

For a grammar is only string derivable from the start symbol. Formally $S \Rightarrow^* \alpha$

The language of a grammar is the set of terminal strings such that the start symbol w derives w

eg G has production $S \rightarrow \epsilon$ and $S \rightarrow 0S1$
 $L(G) = \{0^n 1^n \mid n \geq 0\}$

CONTEXT FREE LANGUAGES

A context-f. language is a language defined by some context-free grammar.

Some languages are not context-free. Because context-f. grammars can count two things, but not three eg $0^n 1^n 2^n \mid n \geq 1$

BNF notation (Backus-Naur form)

Grammars used for programming languages often

- In BNF you usually use a word to describe a variable. For example, statement for the intent is that this variable will generate all the strings that are valid statements for the programming language.

Conventionally these words are put in brackets that tell they are variables

- Terminals are $i + * \text{ while if}$

$::=$ is used for reduction

$|$ is used for lists of productions with some left if

\dots is similar to the , but means one or more, instead

$\langle \text{digit} \rangle ::= 1|2|3|4|5|6|7|8|9$

$\langle \text{unsigned int} \rangle ::= \langle \text{digit} \rangle \dots \langle \text{digit} \rangle$

TRANSLATION we can replace \dots with a new variable A and productions $A \rightarrow \alpha \mid \beta$
 eg $U \rightarrow U|D \quad D \rightarrow 0|1|2|3|4|5|6|7|8|9$

OPTIONAL ELEMENTS ARE ADDED

by surrounding with $[\]$

eg $\langle \text{statement} \rangle ::= \text{if} \langle \text{condition} \rangle \text{ then } \langle \text{statement} \rangle [\text{else } \langle \text{statement} \rangle]$

TRANSLATION replace $[\alpha]$ by a new variable A with productions $A \rightarrow \alpha \mid \epsilon$

$S \rightarrow iCtSA \quad | \text{if } S \text{ then } S \text{ else } S$
 $A \rightarrow ;eS \mid \epsilon$
 C for concatenation
 A is an optional list
 S is a statement
 i is a terminal
 t is a terminal
 e is a terminal

Say $\alpha \Rightarrow_{\text{TM}} \beta$ if w is a string of terminals only $A \rightarrow \beta$ is a production

$S \Rightarrow_{\text{TM}} SS \mid (S) \mid ()$

$S \Rightarrow_{\text{TM}} SS \Rightarrow_{\text{TM}} S(S) \Rightarrow_{\text{TM}} (S)(S) \Rightarrow_{\text{TM}} ((S))$

$S \Rightarrow_{\text{TM}} ((S))$

$S \Rightarrow_{\text{TM}} (((S)))$

PARSE TREES

A parse tree is a graph that shows how some string is derived using some context-free grammar

They let us express the concept of ambiguity in grammars. Unambiguity, that the ability to provide a unique tree structure for every string in its language is vital.



The root is usually labeled as the start symbol. But not always (if there is a root A)

The yield of a parse tree is the string of labels from the left

A CFG grammar is ambiguous if there is a string in the language that is the yield of two or more parse trees.

In a grammar where you can always figure out the production by scanning the given string left-to-right and looking at the next symbol is called LL(1) (Most programming languages)

There are equivalent unambiguous grammar for some language. When this is not possible, we say the language is **inherent ambiguous**.

Use $\{ \}$ to group several different elements. eg $\langle \text{statement list} \rangle ::= \langle \text{statement} \rangle [\{ ; \langle \text{statement} \rangle \}^*]$ optional

For programming languages where you have a list of statements separated by a semicolon

TRANSLATE

To translate to our original notation

- Create a new variable A for $\{ \alpha \}$
 - One production for A : $A \rightarrow \alpha$
 - Use A in place of α

eg $L \rightarrow S[A \dots] \quad A \rightarrow iS$
 A stands for $\{ ;S \}$

$L \rightarrow SB \quad B \rightarrow A \dots \mid \epsilon \quad A \rightarrow ;S$

We introduce a new variable B

$L \rightarrow SB \quad B \rightarrow C \mid \epsilon \quad C \rightarrow AC \mid A \quad A \rightarrow ;S$

We replace $A \dots$ in B with the new variable C

LEFT-MOST and RIGHT-MOST derivation

- Derivations to allow us to replace any of the variables in a string. Leads to many different derivations of the same string
 - By forcing the left-most (or right-most) variable to be replaced, we avoid that "distinction without a difference"

Say $w \Rightarrow_{\text{TM}} w'$ if w is a string of terminals only $A \rightarrow \beta$ is a production

Also $\alpha \Rightarrow_{\text{TM}}^* \beta$ if α becomes β by a sequence of zero or more \Rightarrow_{TM} steps eg $S \Rightarrow_{\text{TM}} SS \mid (S) \mid ()$

$S \Rightarrow_{\text{TM}} SS \Rightarrow_{\text{TM}} (S)S \Rightarrow_{\text{TM}} ((S))S \Rightarrow_{\text{TM}} (((S)))$

NORMAL FORM

Become context-free variables can be badly designed. Analogous to the states of a finite automaton not reachable from the start state. There are productions that take many steps that obviously can be combined.

DISCOVERING ALGORITHMS

Almost all the algorithms used to find local context-free grammars are based on the same principle.

They use induction

- The basis are facts that are obvious
- discover more facts from what it already discovered
- until no more facts can be discovered

Find variables that cannot derive terminals

Basis. If there is the production $A \rightarrow w$ where w has no variables, then A derives terminal string in one step (can be the empty string)

Induction. $A \rightarrow \alpha$ and α consists of terminals and variables previously discovered. Then A derives a terminal string.

Every variable that derives a terminal string will be discovered

Algorithm to eliminate variables that derive nothing

- Find the variables that derive terminal strings
- For the other variables, remove all productions in which they appear in either head or body

Unreachable symbols

Another a terminal or variable deserves to be eliminated is if it cannot be derived from the start symbol

This algorithm is another example of the discovery algorithm

Basis, the start symbol can be reached

Induction, suppose we discovered A , then a production $A \rightarrow \alpha$ says that we can reach all symbols in α

When we can discover no more symbols, then we have all symbols derivable from S . Remove all the other symbols.

Eliminate useless symbols

A symbol is useless if it appears in some derivation of some terminal string from the start symbol

It is useless otherwise

- 1 Eliminate all symbols that derive no terminal string
- 2 Eliminate unreachable symbols

ϵ productions

$A \rightarrow \epsilon$, they can be eliminated from a context-free grammar.

The only thing is that we can no longer derive the empty string.

Theorem If L is a CFL, then $L - \{\epsilon\}$ has a grammar with no ϵ -productions

Nullable symbols

We discover nullable symbols, which are variables A such that $A \Rightarrow^* \epsilon$

Basis, if $A \rightarrow \epsilon$, it is nullable

Induction, if $A \rightarrow \alpha$ and all symbols in α are nullable, then A is nullable

Eliminate ϵ productions

Turn each production $A \rightarrow x_1 \dots x_n$ into a family of productions

For each set of nullable x_i , we delete the nullable production and make a new production (except in the case all productions are nullable)

Unit productions

where the body has a single variable.

They can be eliminated as well

If $A \Rightarrow^* B$ and $B \Rightarrow^* \alpha$ is a non-unit production. Then add the production $A \rightarrow \alpha$

Then drop all unit productions

The algorithm finds all pairs (A, B)

Basis, (A, A)

Induction, if we have found (A, B) , and $B \rightarrow C$ is a unit production, add (A, C)

Cleaning up a grammar

If L is a CFL, then there is a grammar for $L - \{\epsilon\}$ such that
- No useless symbols - No ϵ -productions - No unit productions

Chomsky normal form

Has only: $A \rightarrow BC$ If L is a CFL, then $L - \{\epsilon\}$ has a CFG in CNF
 $A \rightarrow a$

PUSHDOWN AUTOMATA

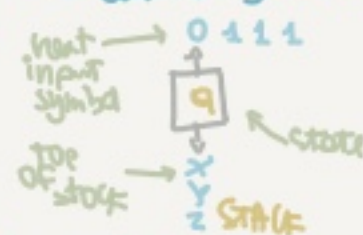
Equivalent in power of context-free grammars

Only the nondeterministic version defines all context-free languages

Deterministic PDAs are used to model parsers

Think of an ϵ -NFA with a stack to store symbols. Moves are determined by

- Current state
- Current symbol (or ϵ)
- Current symbol on stack



A non-det. PDA can have a choice of next moves

- Change of stack
- Modify the stack by zero or more symbols with push or pops

A PDA is described by $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$

stack alphabet

start stack symbol

Conventions

- a, b, c, \dots are input symbols
- \dots, X, Y, Z are stack symbols
- \dots, w, x, y, z are strings of input
- d, β, \dots are strings of stack symbols

Transition function

- Takes three arguments $\delta(q, a, z)$

This is a set of zero or more actions, each consist of a new state p and a string α of stack symbols with which to replace the top symbol z ...

- IF an action of a PDA is $\delta(q, a, z)$ and contain (p, a)
- ... Then it can
- Change state to p
 - Remove a from the top of the stack
 - Replace z on the top of stack by b or

Instantaneous description

An instantaneous description, or ID is a triple (q, w, d) where:

- q is the current state
- the remaining input w
- d is the stack contents (top of stack is the leftmost symbol of d)

There an analogy between derivations in a grammar and sequences of ID's for a PDA

We say ID I can become ID J in one move like that $I \vdash J$

Formally $(q, a, w, x, d) \vdash (p, w, \beta, d)$ if $\delta(q, a, x)$ contains (p, β)

top character a, w is the rest
 x is top symbol
 d is everything below

We can extend \vdash to \vdash^* , meaning zero or more moves

Language of a PDA

IF P is a PDA, then $L(P)$ is the set of strings such that $(q_0, w, z_0) \vdash^* (F, \epsilon, d)$

Final state F anything left on the stack
 input consumed

Another language is defined by the set of strings that make empty its stack

This language is called $N(P)$ for a PDA P

Formally $(q_0, w, z_0) \vdash^* (q, \epsilon, \epsilon)$ q can be non-final

IF $L = L(P)$, there is a PDA P' such that $L = N(P')$

Deterministic PDA

There must be at most one rule for only state q , input a , and stack symbol x

There must not be a choice between using ϵ or real input

Formally $\delta(q, a, x)$ and $\delta(q, \epsilon, x)$ cannot both be non-empty

PUMPING LEMMA FOR CFG

We can always find two pieces of any long string and after repeating them the same number of times, we will get a string in the same language

Statement

For every context-free L there is an integer n , such that for every string z in L of length $\geq n$ there exists $z = uvwxy$ such that

- $|vwx| \leq n$
- $|vx| > 0$
- For all $i \geq 0, uv^iwx^iy$ is in L

eg

$\{0^i 1 0^i \mid i \geq 1\}$ is a CFL

$\{0^i 1 0^i 1 0^i \mid i \geq 1\}$ is not

- We can't match two pairs or three counts of a group

Proof

- Suppose it is a CFL
- Let n be L 's pumping lemma's constant

adversary picks

consider $z = 0^n 1 0^n$

We can write $z = uvwxy$ where $|vwx| \leq n$ and $|vx| \geq 1$

adversary picks

Case 1, vx has no 0's

- Then at least one of them is a 1, which no string in L does

Case 2, vx has at least one 0

- vwx is too short ($\leq n$) to extend all three blocks of 0's

- Thus $uvxy$ has at least one block of n 0's and at least one block with fewer than n 0's

This part of automata theory lets us show that certain tasks are impossible, or intractable

UNDECIDABILITY

Integers, strings and other, are all represented as strings of bits

KEY POINT Strings that are programs, are also strings of bits

Everything is an integer

But since binary strings can have leading 0's, there is not a unique representation

eg 101, 0101, 00101

To make the correspondence one to one, we put a one in front of any binary string, and then treat it as an integer

101 → 1101 0101 → 10101 00101 → 100101

A set is finite if it has a integer that is the count of the members on the set.

Formally, is one for which it is impossible to find a one to one correspondence between the members of the set and a subset of that set

... And an infinite set is one for which there are to one correspondence

A countable set has a one to one correspondence with positive integers (All countable sets are infinite)

Binary strings are countable. Every program can be represented with a unique integer.

We call the one to one correspondence between a countable set and a positive integer an enumeration of the set

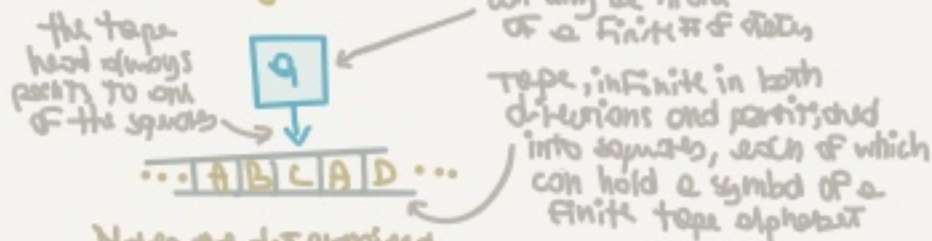
TURING-MACHINE THEORY

The languages over $\{0,1\}$ are NOT countable

The purpose of this is to prove certain languages have no membership algorithm

- start with a language about Turing machines themselves
- reductions are used to prove more complex questions undecidable

Turing machine



Notes are determined by the state and the symbol under the tape head
In the move it can change state, write a new tape symbol over the old one in the square it's scanning, and move the head one square left or right

They are as powerful as computer

... arguably even more because of the infinite tape

Formally, a TM is described by

- A finite set of states, Q
 - An input alphabet, Σ .
 - A tape alphabet, Γ , typically contains Σ
 - A transition function, δ
 - A start state, $q_0 \in Q$
 - A blank symbol, B , in $\Gamma - \Sigma$ typically
- All the tape except for the input is blank initially
- A set of final states

ab... are input symbols
...xyz are tape symbols
...wxyz are strings of input
aβ... are strings of tape symbols

notion of everything we can compute → They both define the same class of languages and their subset, the **RECURSIVE LANGUAGES**

An algorithm is a TM accepting by final state, that is guaranteed to halt whether or not it accepts
If $L(M)$ for a TM is an algorithm, then L is a recursive language

The transition function

$\delta(q, z)$ for state q and tape symbol z is either undefined, or a triple (p, y, D) where p is a state, y is a new tape symbol, and D is a direction L or R

Instantaneous description

Initially a TM has a tape consisting of a string of input symbols, surrounded by an infinity of blanks in both directions

The TM begins in its start state, with the head at the left-most input symbol

An ID has the form $\alpha q \beta$ where $\alpha \beta$ includes the tape between the leftmost and rightmost non-blanks

The state q is immediately to the left of the symbol the head is now scanning

We use \vdash to represent a move and \vdash^* to represent any number of moves

Moves, formally

If $\delta(q, z) = (p, y, R)$ then $\alpha q z \beta \vdash \alpha y p \beta$
↳ if z is blank, then $\alpha q \vdash \alpha y p$

If $\delta(q, z) = (p, y, L)$ then for any x , $\alpha x q z \beta \vdash \alpha p x y \beta$
- in addition $q z \beta \vdash p B y \beta$ $\alpha p x y \beta$

blank here move the blank

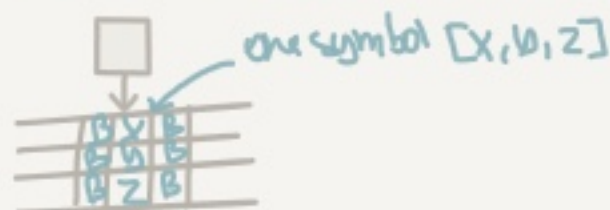
Languages of a TM

$L(M) = \{w \mid q_0 w \vdash^* I \text{ where } I \text{ is an id with a final state}\}$

Or by halting, $H(M) = \{w \mid q_0 w \vdash^* I \text{ and there is no move from ID } I\}$

TRICK: MULTIPLE TRACKS

- Think of tape symbols as vectors with k components, each chosen from a finite alphabet
- Makes the tape appear to have k tracks
- Input symbols have the blank in all components but one



TRICK: MARKING

A good use of multitasking is to use one track for data, and another track for marks

In the mark track, almost all squares have blank value, but one or more have special symbols (marks) that indicate a place on the tape

TRICK: CACHING THE STATE

By thinking the state as a vector, with each component from some finite alphabet, we can use the first component to control operations, what we normally think as a state, and other components to hold data from a finite alphabet (typically bits of tape symbols)

eg this TM copies input w infinitely (fairly simple)

control states

- q : mark the current position and remember the input symbol seen
- p : Run right, remembering the symbol, looking for a blank. Deposit symbol when found it
- r : Run left looking for mark left by q , when found it remove the mark and enter state q

A state is a vector $[x, y]$. x is p, q or r and y is the cache with alphabet $0, 1, B$

...

Tape symbols have the form $[U, V]$. The first component is the marking track, so it is either x (a mark) or B . The other component is the data track (input symbol) and it can be $0, 1$ or B

$$\delta([q, B], [B, a]) = ([p, a], [X, a], R)$$

- In state q , copy the input symbol under the head (a) into the cache
- Mark the position read with X
- Move the head right (R)

a and b represent either 0 or 1

$$\delta([p, a], [B, b]) = ([p, a], [B, b], R)$$

In state p , if the current tape symbol does not have a blank, stay in that state and move right

$$\delta([p, a], [B, B]) = ([r, B], [B, a], L)$$

- when you find a B , place a (in the cache) in the data track
- go to state r and move left (L)

$$\delta([r, B], [B, a]) = ([r, B], [B, a], L)$$

- In r , move to the left looking for the mark

$$\delta([r, B], [X, a]) = ([q, B], [B, a], R)$$

When the mark is found we

- remove the mark
- go to state q
- move to the right

we're in state q , the whole cycle repeats

...

DECIDABILITY

To enumerate TMs, we encode them as binary strings

- our alphabet is $\{0,1\}$
- assign integer codes to the components of the TM
 - states: $q_1(\text{start}), q_2(\text{final}), q_3, \dots$
 - symbols: $x_1(0), x_2(1), x_3(\text{blank}), x_4, \dots$
 - directions: $D_1(L), D_2(R)$

Suppose $\delta(q_i, x_j) = (q_k, x_l, D_m)$

- It is represented as $0^i 1 0^j 1 0^k 1 0^m$

All integers are > 0 , so there are no consecutive 1s
of 0s, separated by 1s

Represent a TM by concatenating the encoding for each move, separated by 1 as punctuation
that is $\text{code}_1 11 \text{code}_2 11 \text{code}_3 11 \dots$

Once we have the encoding for our TM, we can convert that binary string to a unique integer that way

- put a one in front of the binary string and treat the result as a binary integer

Thus it makes sense to talk about the i -th binary string and about the i -th TM

NOTE
Some binary strings represent flawed TMs (as in the wrong place for example). We assume that those machines represent the empty language

Table of TMs related to the strings they accept

TM	1	2	3	4	5	6	7
1							
2							
3							
4							
5							
6							
7							

when 0, the TM does not accept the string. It does if 1

Whenever we have a table like this, we can diagonalize it

Formally $D = a_{i,i}$ where $i = 0$ if the (i,i) table entry is 1 and vice-versa

D is not a row in the matrix, therefore it does not represent a language accepted by any Turing machine

- Suppose it were the j -th row
- D disagrees with the j -th row at the j -th column
- Thus D is not a row

Consider the diagonalization language

$L_D = \{w \mid w \text{ is the } i\text{-th string and the } i\text{-th TM does not accept } w\}$

Since L_D has no TM, it is not recursively enumerable

Problems

Informally, a problem is a yes/no question about an infinite set of possible instances

eg. Does graph G have a Hamilton cycle (cycle that touches each node exactly once)?

- Each undirected graph is an instance of the "Hamilton-cycle problem"



Here the answer is yes, there is a Hamilton cycle

- Formally, a problem is a language over some alphabet Σ

- Each string in Σ^* can be viewed as an instance of the problem once we decide on encoding

- The string is in the language if the answer to this instance of the problem is "yes".

NOTE
For flawed encodings, we assume the answer is "no"

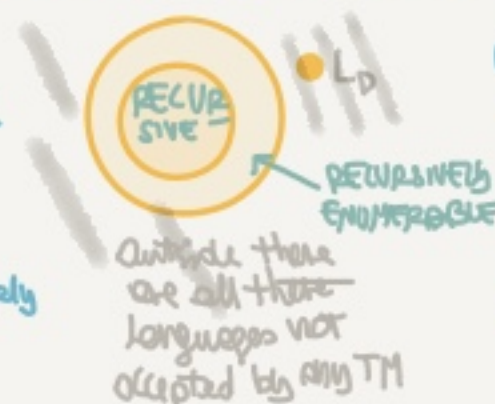
We can think of the language L_D as a problem
Does this TM not accept its own code?

Decidable problems

A problem is decidable if there is an algorithm to answer it

Then a decidable problem is a recursive language

Otherwise a problem is undecidable



Undecidable problem are either in the outer ring (RE languages but not recursive) or in the outer space.

Real world undecidable problems could be

- Can a particular line of code ever be executed?
- Is a given context-free grammar ambiguous?
- Do two CFB's generate the same language?

A recursively enumerable, but not recursive, language is the language L_U of a universal Turing machine

An UTM takes as input a binary string w and a TM M . w is accepted if and only if M accepts w .

eg. A Java virtual machine takes a coded Java program and an input for that program and then it executes that program on the input. A JVM is a more general UTM as it does not create a single accept output.

Designing the UTM

Inputs are of the form:

Code for $M + 111 + w$

a valid TM has never 111, so we can avoid ambiguities

Has several tapes

- Tape 1 holds the input $M111w$
- Tape 2 is used to simulate the tape of M
- Tape 3 holds the state of M

Step 1 the UTM checks that M is a valid code for a TM

If M is not valid, its language is empty, so the UTM immediately halts without accepting

Step 2 The UTM examines M to see how many of its own tape squares it needs to represent one symbol of M

Step 3 Initialise Tape 2 to represent the tape of M with input w , and initialise tape 3 to hold the start state

Step 4 Simulate M .

- Look for a move on tape 1 that matches the state on tape 3 and the tape symbol under the head on tape 2
- If found, change the tape symbol and move the head marker on tape 2 and change the state on tape 3
- If M accepts, the UTM also accepts

INTRACTABLE PROBLEMS

Decidable problems, but they take at least exponential time as a function on their input size

In practice, those problems can only be solved for small instances

An unsolvable problem in practice, has to run in less than exponential time, polynomial time, for some low degree polynomial

Time-Bounded TM's

A TM that, given an input of size n , always halts within $T(n)$ moves is said to be $T(n)$ -time bounded

- TM can be multitape
- Sometimes, it can be non-deterministic

The class P

- If a DTM M is $T(n)$ -time bounded for some polynomial $T(n)$, then we say that M is polynomial time
- $L(M)$ is said to be in the class P
- When we talk of P, doesn't matter whether we mean "by a computer" or "by a TM" (If a computer takes $O(T(n))$ for a given algorithm, a TM can simulate it in at most $O(T^2(n))$ steps)

Running times between polynomials

P only requires that the language be accepted by some TM whose running time is bounded above by some polynomial ($O(\log n)$ or $O(n \log n)$ are included)

KNAPSACK, seems to be in P but isn't

- Given positive integers i_1, \dots, i_n , can we divide them in two sets with equal sums
- At first glance, we can solve it by using dynamic programming. Maintain a table of all the differences by partitioning the first j

Although that runs in no more than the square of the sums of the integers, it doesn't tell us that knapsack is in P.

The problem is that we cannot define the input size as the sum of the integers in the input. The input size is always the # of cells that need to write the input on a TM tape.

- Suppose we have n integers, each of which is 2^n
- We can write integers in binary, so the input takes $O(n^2)$
- But the table requires space $O(n2^n)$
- All n tables in time $O(n2^n)$. By using n as the input size, it requires $O(n^{2 \sqrt{n}})$

The class NP

The running time of a nondeterministic TM is the maximum number of steps taken along any branch (makes any sequence of choices)

If that time bound is polynomial, the NTM is said to be polynomial-time bounded

And its language is said to be in the class NP

Knapsack is NP:

First, uses its non-determinism to guess a partition of the integers into subsets

...then sum the subsets and compare, say yes if two partitions yield two equal sums.

P versus NP

Are P and NP the same class of languages? That is, can any problem that is solved by a NTM in polynomial time also be solved by some det. TM in polynomial time even if the degree is higher?

Question posed by Steve Cook in 1970. After all, NFA can be simulated by a DFA, and NTM's can simulate det. ones. But the problem is still very very difficult.

- There are thousands of NP problems for which no algorithm in P has been found. Neither there is a proof that these problems are not in P