

from neural networks to the structure of language: a physicist's perspective

Eric DeGiuli

Institut de Physique Théorique Philippe Meyer
École Normale Supérieure, Paris

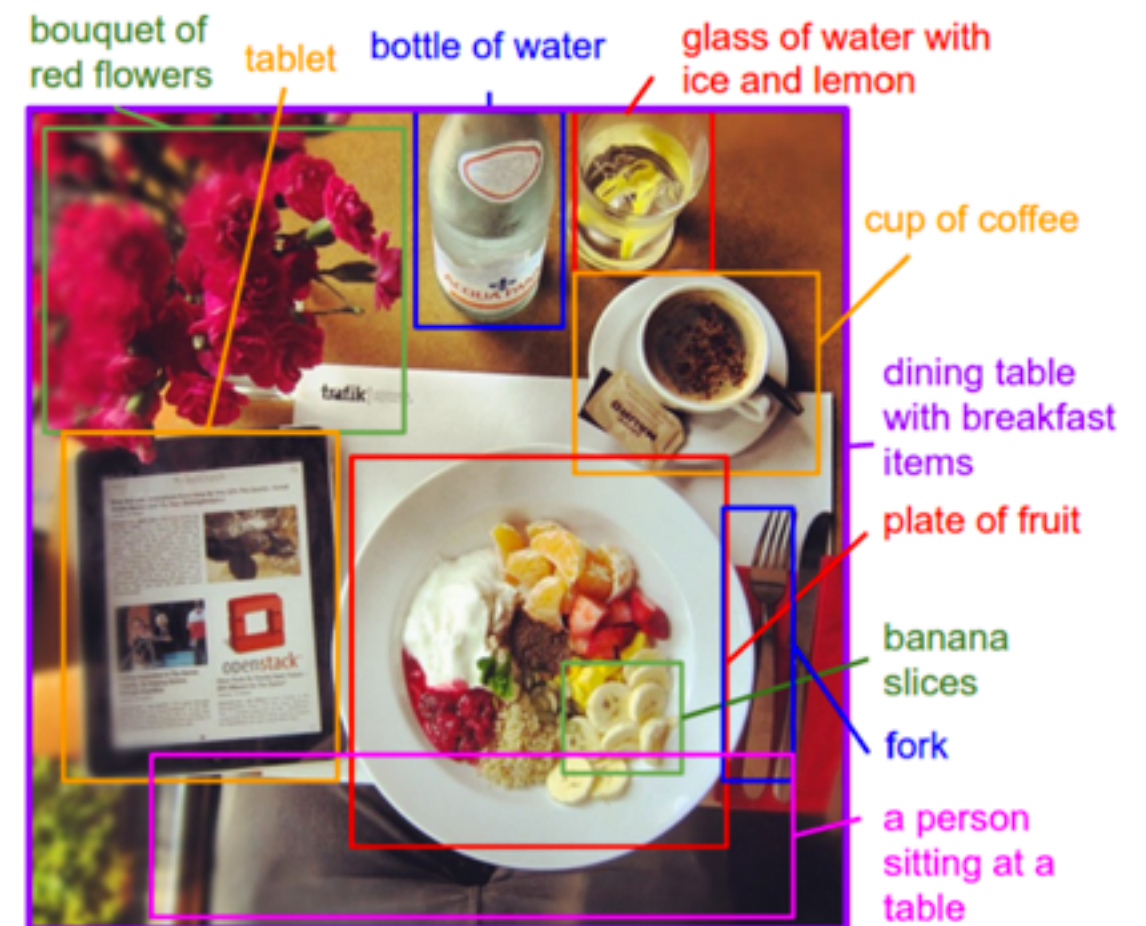


successes of machine learning

speech recognition



image recognition

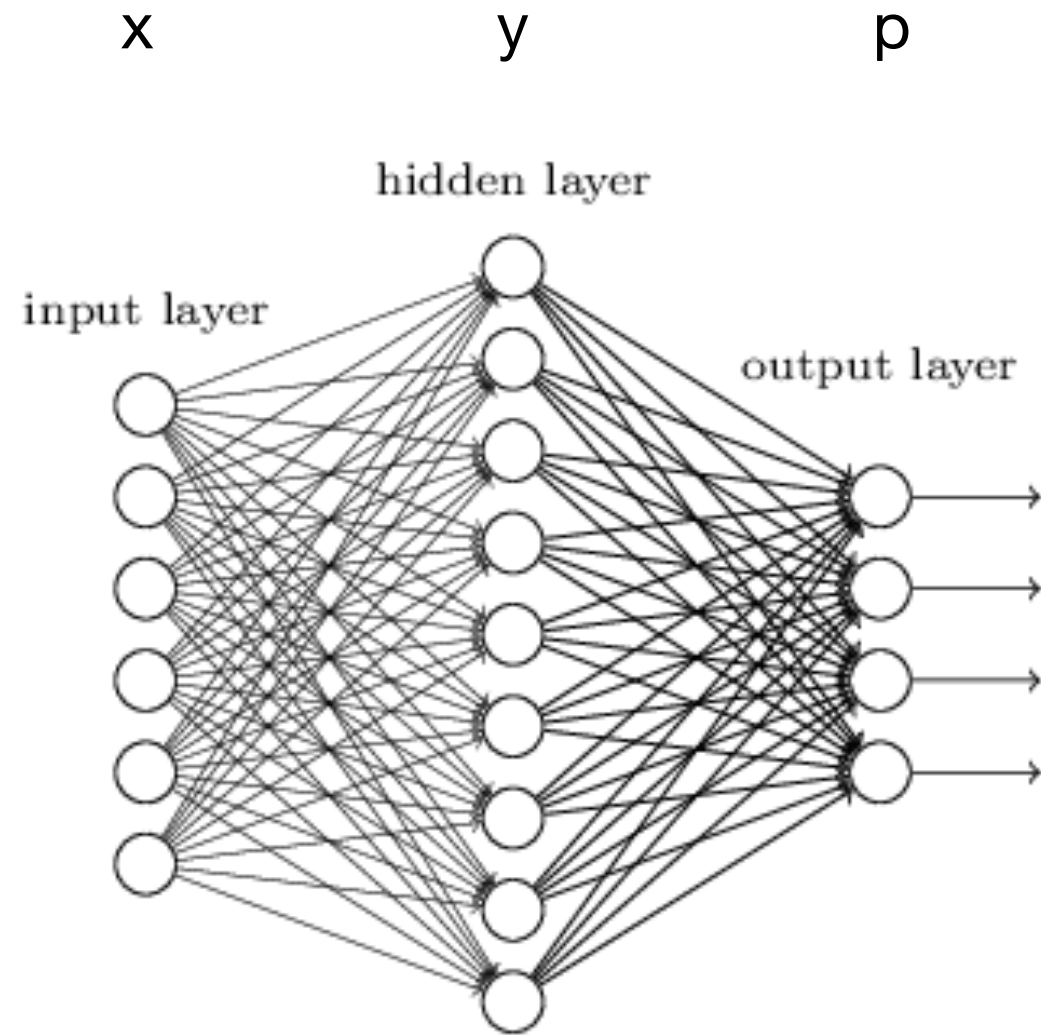


Example output of the model

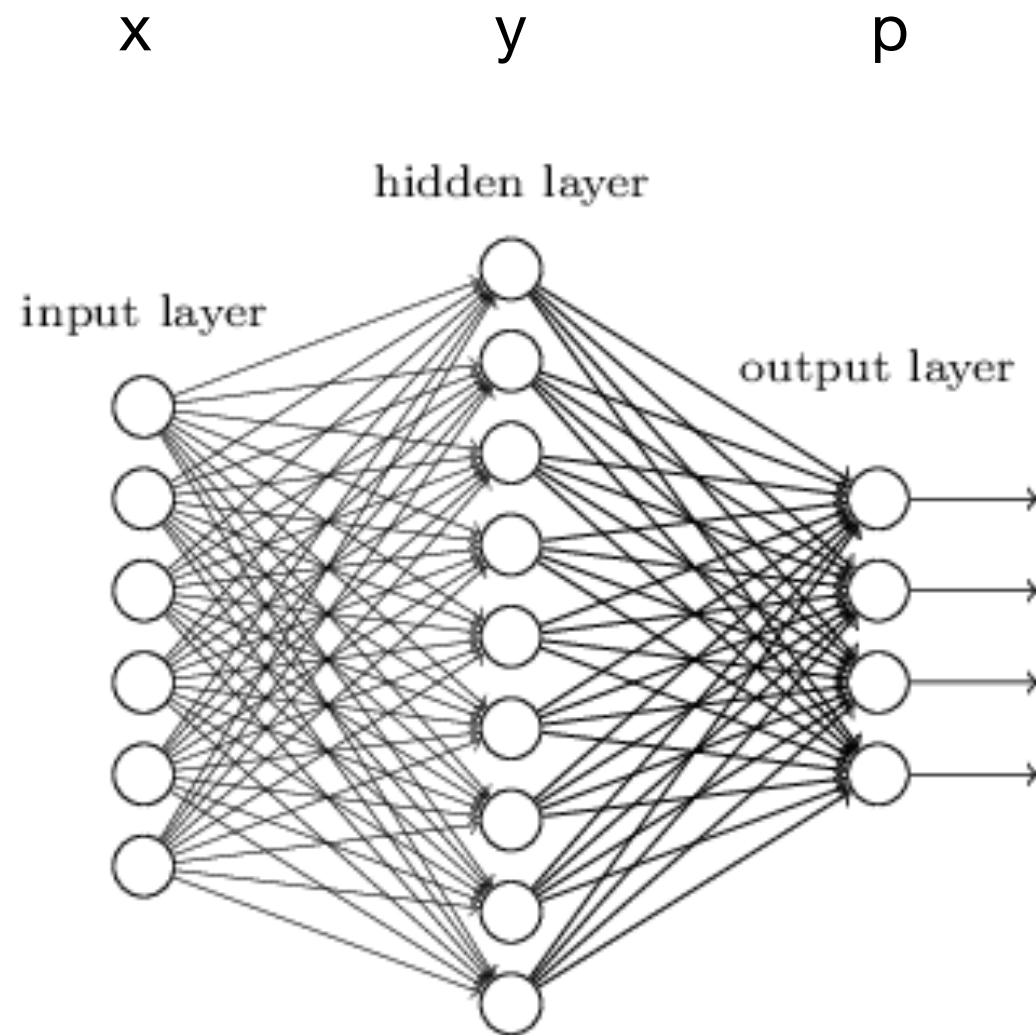
all using deep neural networks

Karpathy & Li. *Proc. IEEE CVPR* 2015

neural network architecture



neural network architecture



x = input

y = hidden variables
 $= f(Ax + b)$

A = parameter matrix

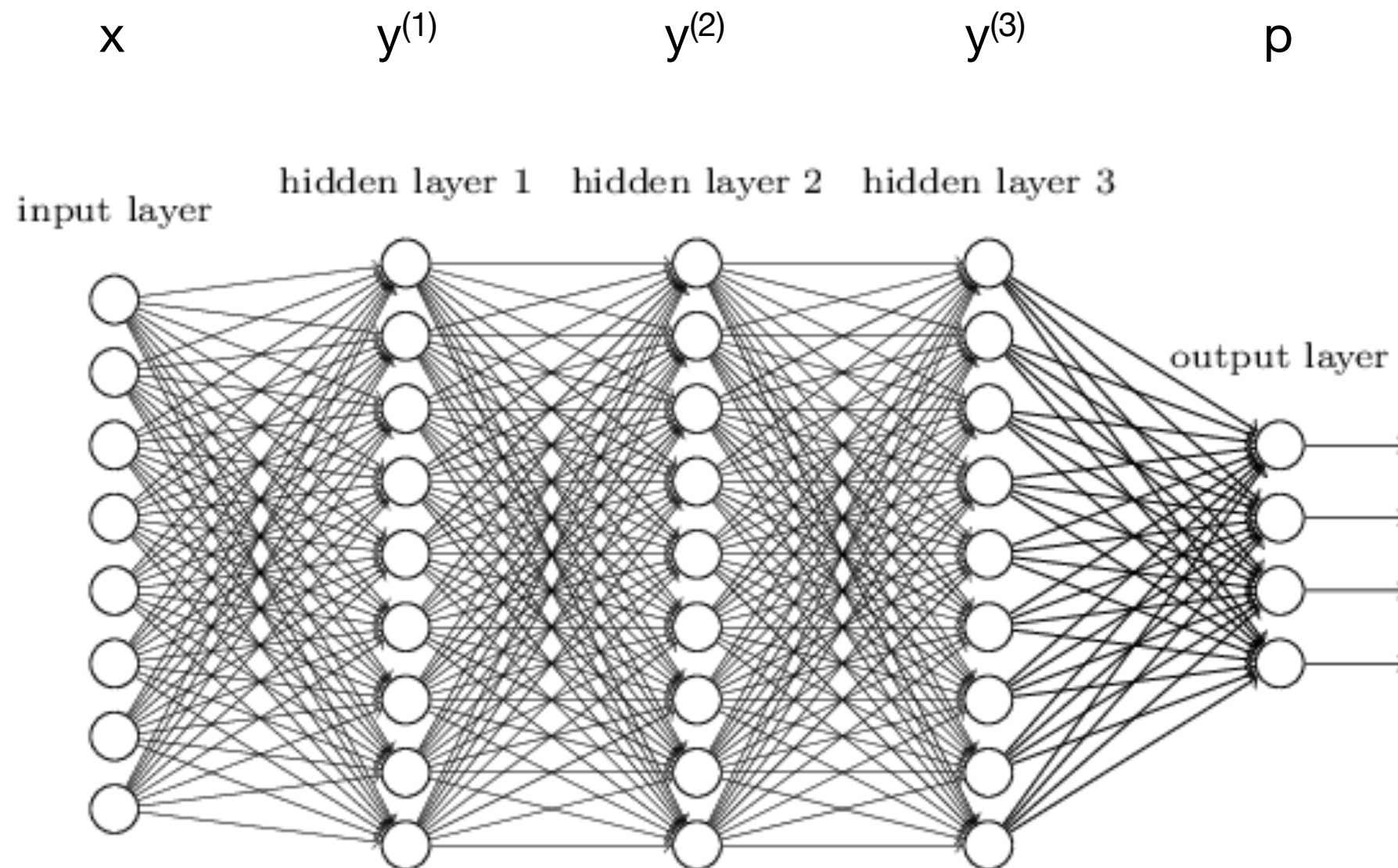
b = offset parameter vector

f = component-wise activation
function

p = probabilities
 $= g(y)$

goal of learning: determine optimal A , b

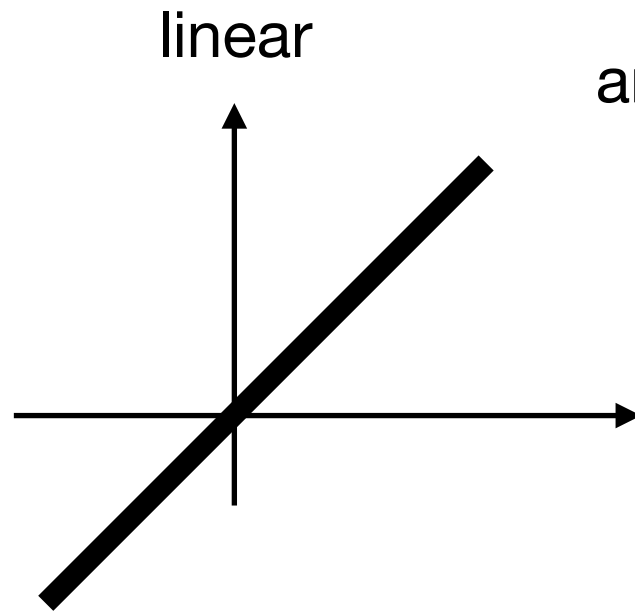
neural network architecture



deep neural network = network with several hidden layers

activation

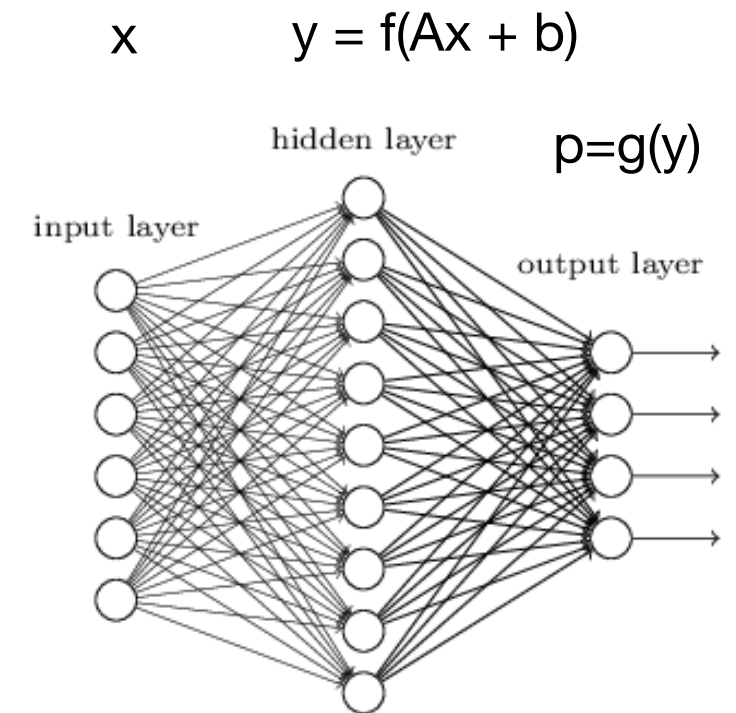
how to choose f ?
consider dp vs dx



any change in input leads to change in y

all elements contribute to dp

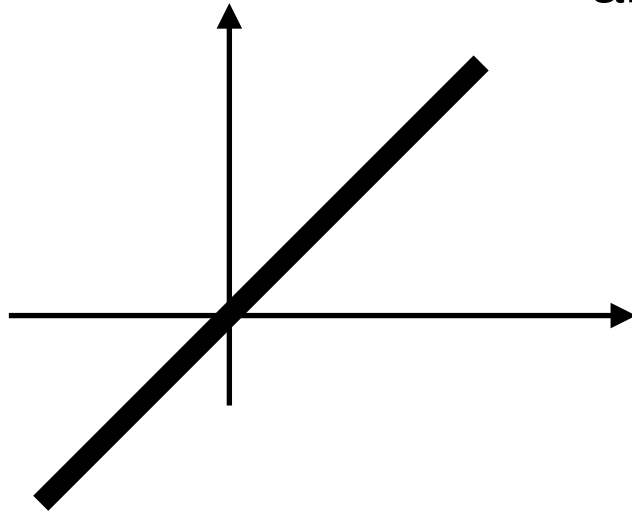
referendum machine



activation

how to choose f ?
consider dp vs dx

linear

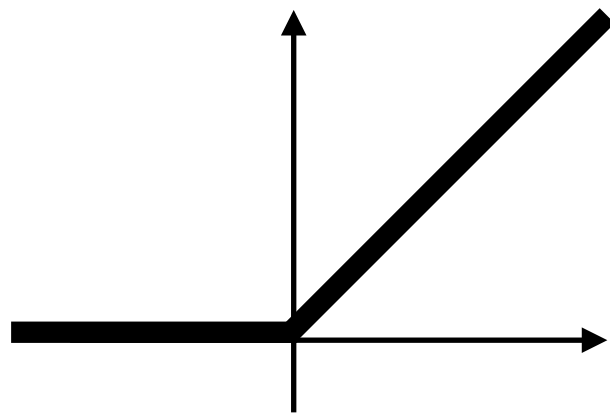


any change in input leads to change in y

all elements contribute to dp

referendum machine

nonlinear



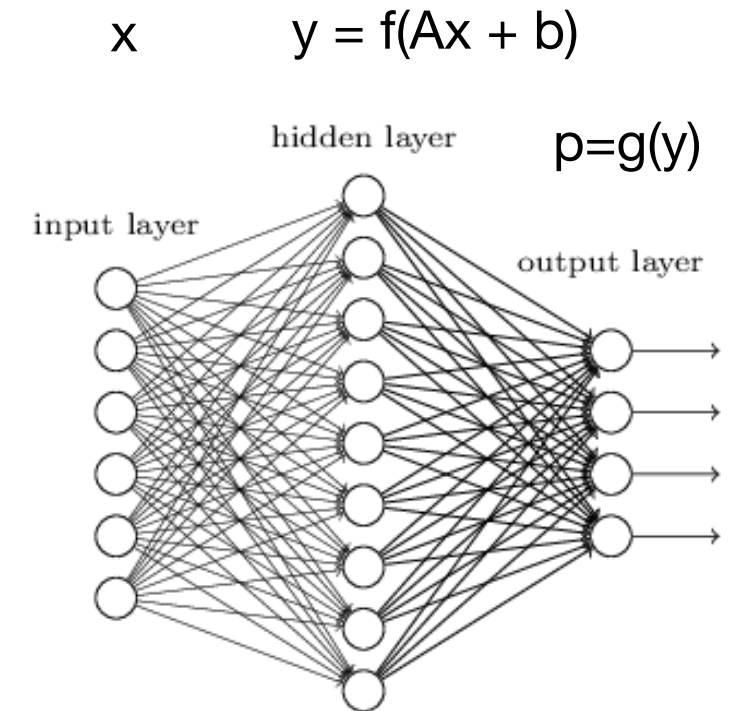
common choice: clamp¹

$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

no change in output for $x < 0$

elements can be indifferent

expert machine



¹ 'ReLU' = rectified linear

composition

- essential to have nonlinear activation
- saturation in activation \Rightarrow elements can act compositionally
(expert elements rather than jack-of-all-trades)
- allows approximate factorization of data space \Rightarrow fewer parameter DOF
- theory? (see Tubiana Monasson PRL 2016)

1st principle of successful machine learning: composition

composition

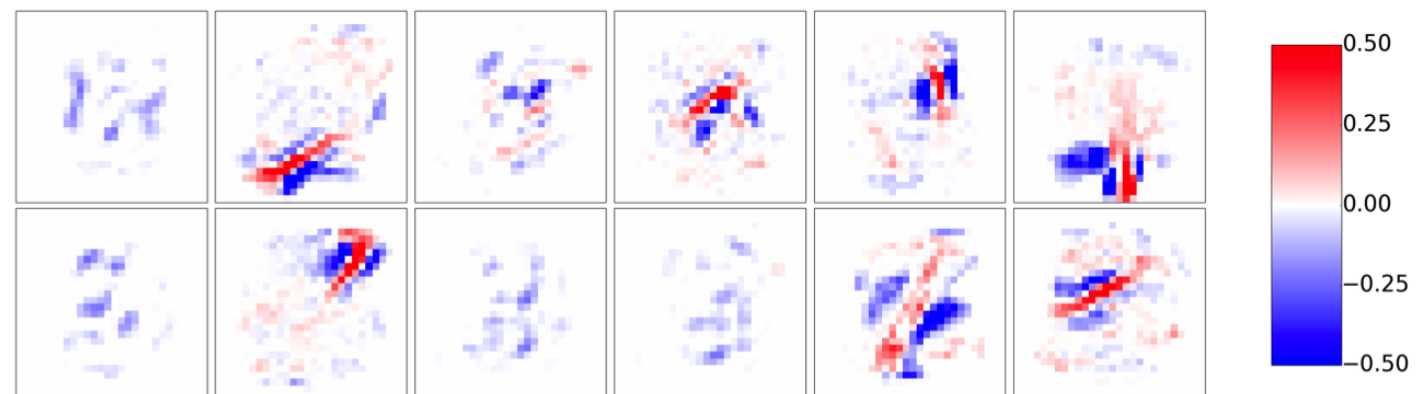
- essential to have nonlinear activation
- saturation in activation \Rightarrow elements can act compositionally
(expert elements rather than jack-of-all-trades)
- allows approximate factorization of data space \Rightarrow fewer parameter DOF
- theory? (see Tubiana Monasson PRL 2016)

1st principle of successful machine learning: composition

e.g. handwritten digits



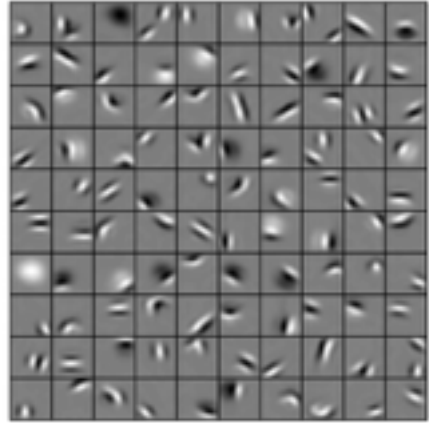
hidden units ~ elementary strokes



Tubiana Monasson 2016

hierarchy

what is role of deep structure?



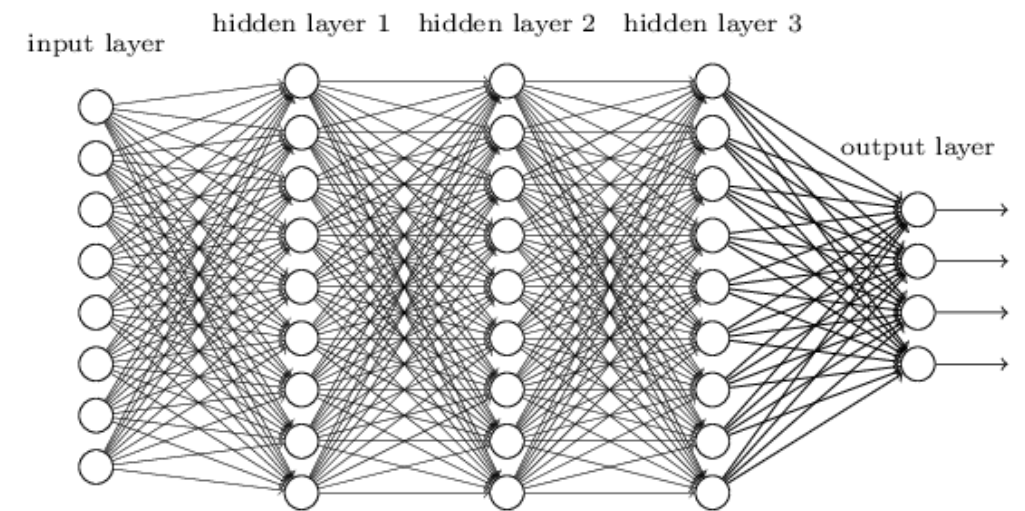
$y^{(1)}$



$y^{(2)}$

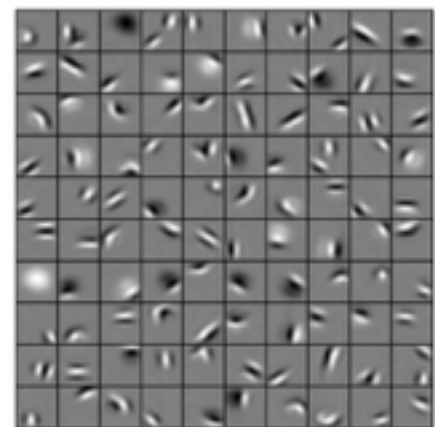


$y^{(3)}$



Lee, H, et al. *Comm. ACM* 54.10 (2011): 95-103.

hierarchy



$y^{(1)}$

what is role of deep structure?

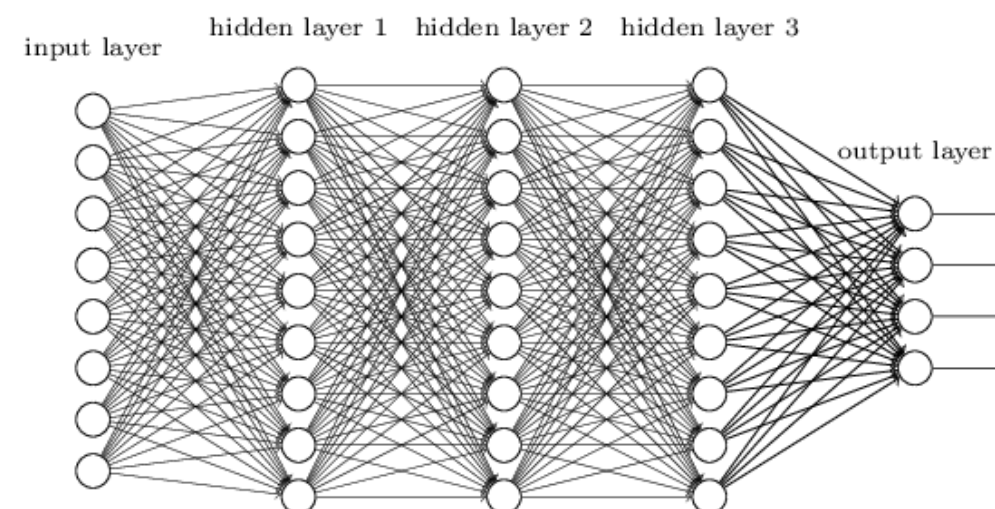


$y^{(2)}$

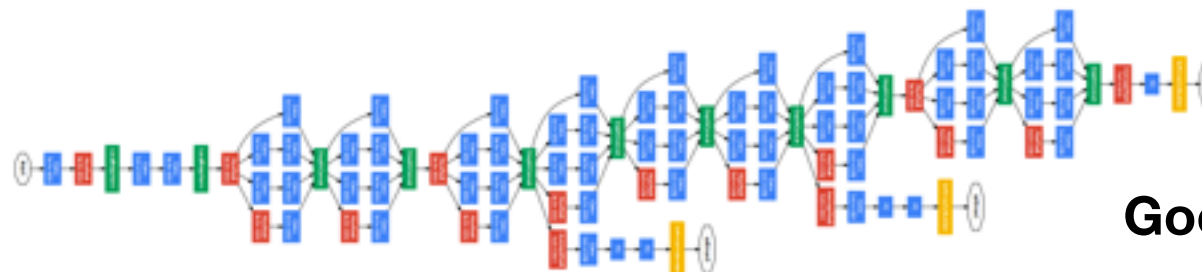
$y^{(3)}$



Lee, H, et al. *Comm. ACM* 54.10 (2011): 95-103.



- deep structure \Rightarrow hierarchical features
- can represent functions with exponentially fewer parameters (Lin, Tegmark, Rolnick J.Stat.Phys 2017)
- empirically, deeper = better



GoogLeNet

Szegedy et al. *Cvpr*, 2015.

2nd principle of successful machine learning: hierarchy

neural network paradigm

architecture: composition & hierarchy

training?

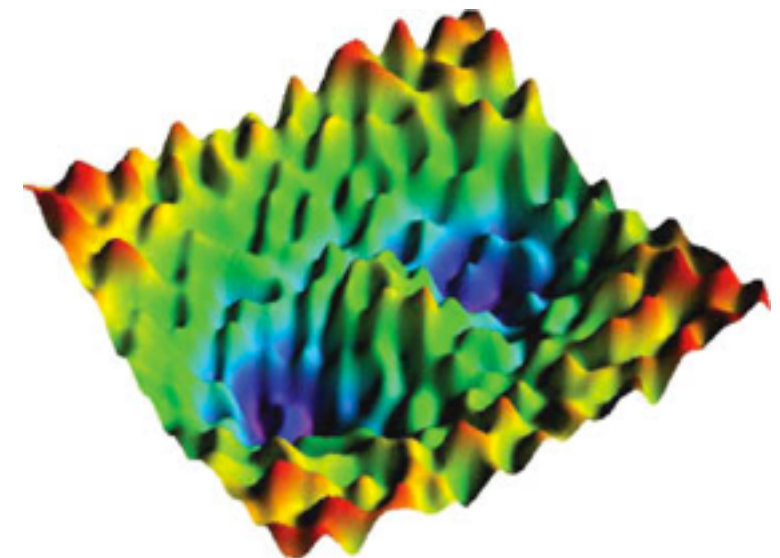
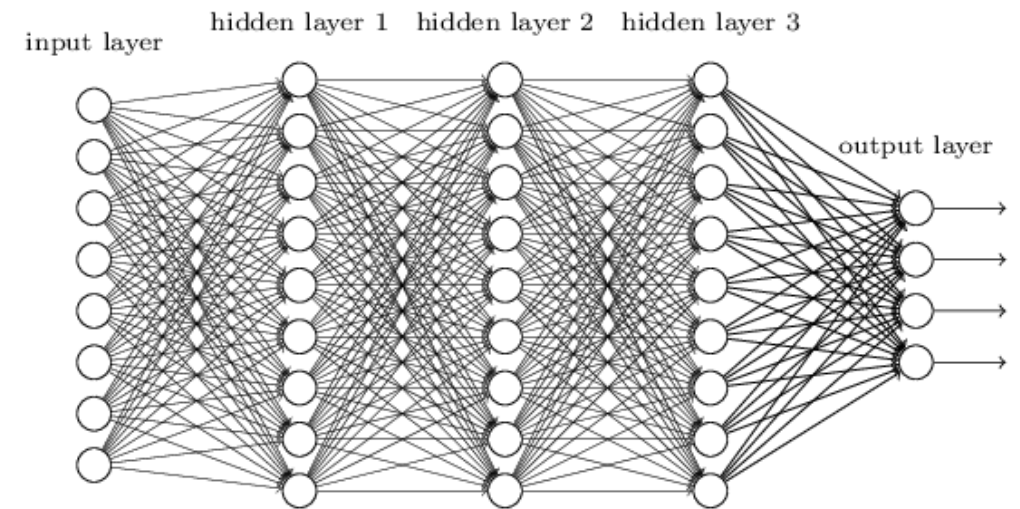
consider 'supervised' learning

have training data pairs $(x_{\text{data}}, p_{\text{data}})$ known exactly

define 'energy' $E(A,b) = \sum (p(x_{\text{data}}) - p_{\text{data}})^2$

- minimizing E is a disordered physics problem (disorder = fixed training data)
- do gradient descent
- phenomenology ~ classic glassy systems

what are the principles for learning?



limitations

neural networks now (2018) are still very far from human intelligence!

e.g. Winograd challenge

1. The city councilmen refused the demonstrators a permit because they feared violence.
2. The city councilmen refused the demonstrators a permit because they advocated violence.

Give 1. Ask `who feared violence?'

Give 2. Ask `who advocated violence?'

humans: > 90%

state-of-the-art (2016): 58%

what is the structure that makes these questions easy for us?

a personal goal:

teach a machine to read & understand a book

why?



© Global Robots Limited

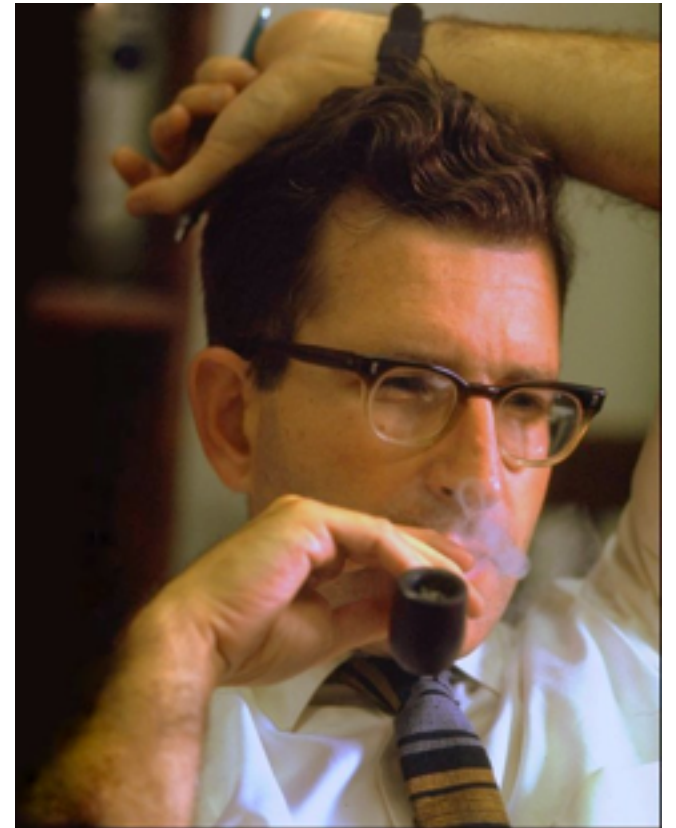
	2013	2014	2015	2016	
...	720,968	887,502	809,128	1,140,078	# articles added to PubMed each year

rigidity of language

1. Is John the man who is tall?
2. *Is John is the man who tall?

rigidity of language

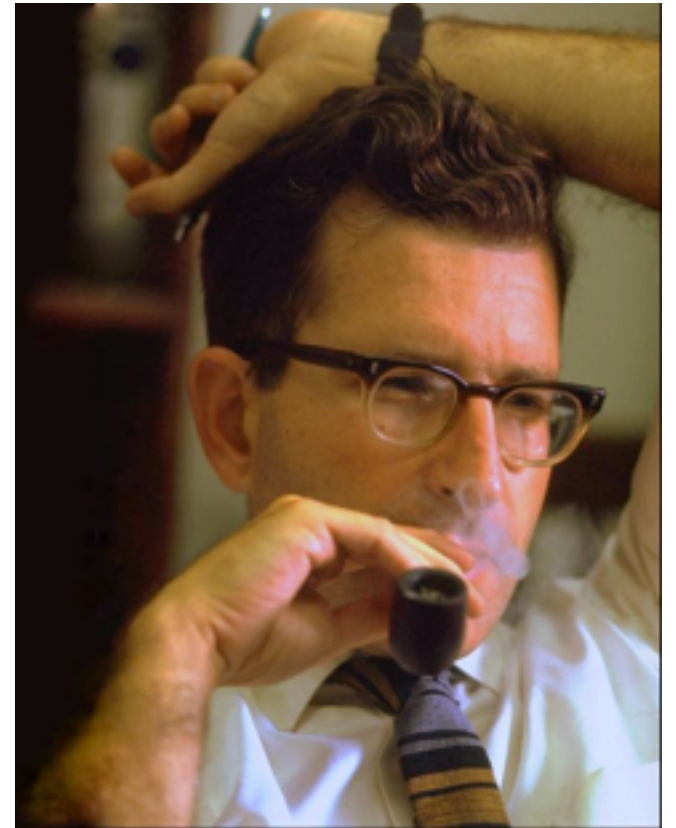
1. Is John the man who is tall?
2. *Is John is the man who tall?
3. Colorless green ideas sleep furiously.
4. *Furiously sleep ideas green colorless.



rigidity of language

1. Is John the man who is tall?
2. *Is John is the man who tall?
3. Colorless green ideas sleep furiously.
4. *Furiously sleep ideas green colorless.

syntax = logical structure
semantics = `meaning' = connection to `truth'



formal grammars

(Pāṇini 400BC, Chomsky, Backus 1950s)

grammar¹ = set of string rewriting rules

A,B,C,.... hidden² symbols

a,b,c,.... observable³ symbols

begin with start symbol, S

repeatedly apply rules until string of
observables

¹ grammar = 'generative grammar' ² 'nonterminal' ³ 'terminal'

formal grammars

(Pāṇini 400BC, Chomsky, Backus 1950s)

grammar¹ = set of string rewriting rules

A,B,C,.... hidden² symbols

a,b,c,.... observable³ symbols

e.g. $S \rightarrow SS$
 $S \rightarrow aSb$
 $S \rightarrow ab$

begin with start symbol, S

repeatedly apply rules until string of
observables

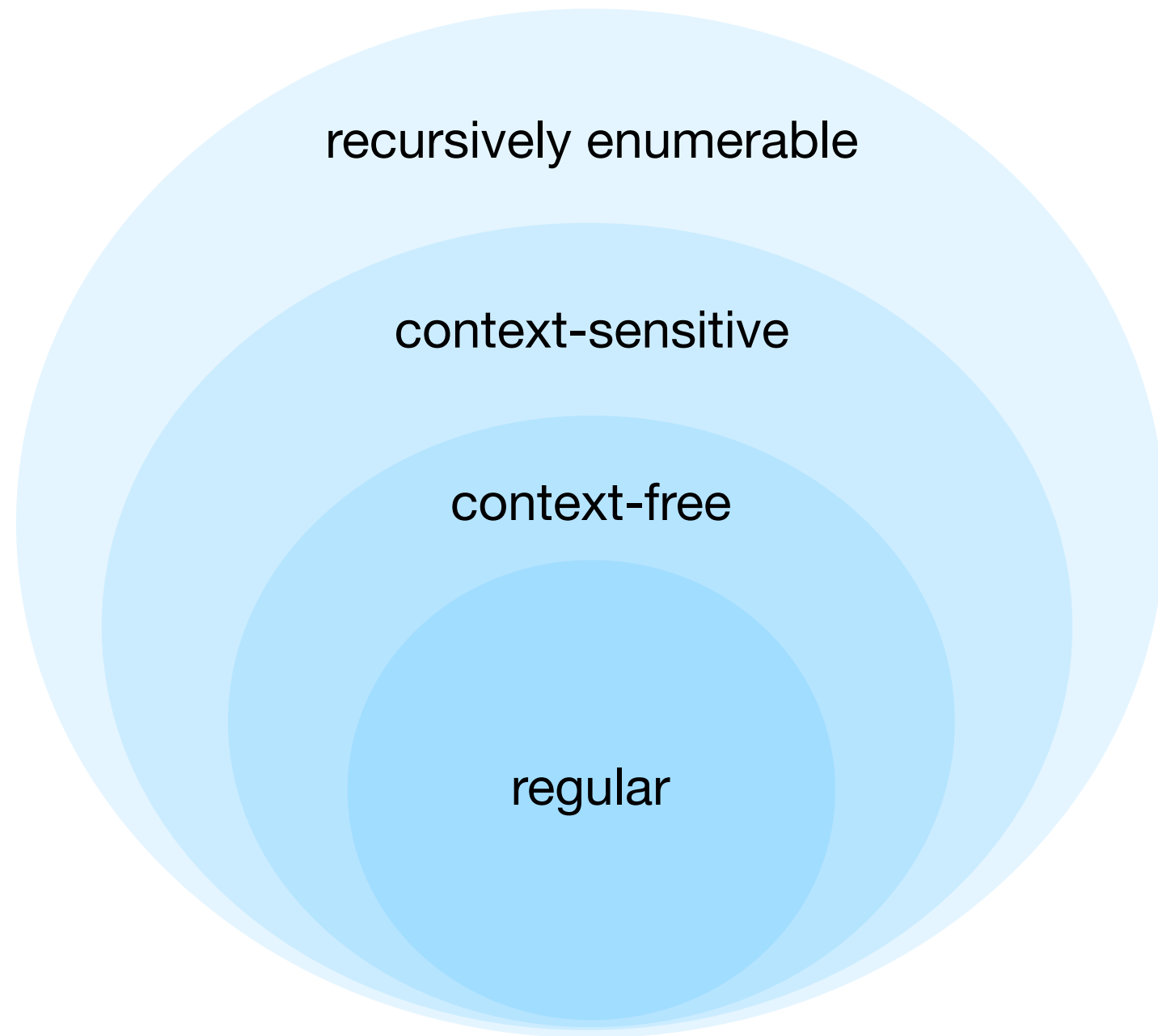
$S \rightarrow SS \rightarrow aSbS \rightarrow aabbS$
 $\rightarrow aabbab$

equivalent to (()) ()

language = set of observable strings

¹ grammar = 'generative grammar' ² 'nonterminal' ³ 'terminal'

Chomsky hierarchy (1950's)

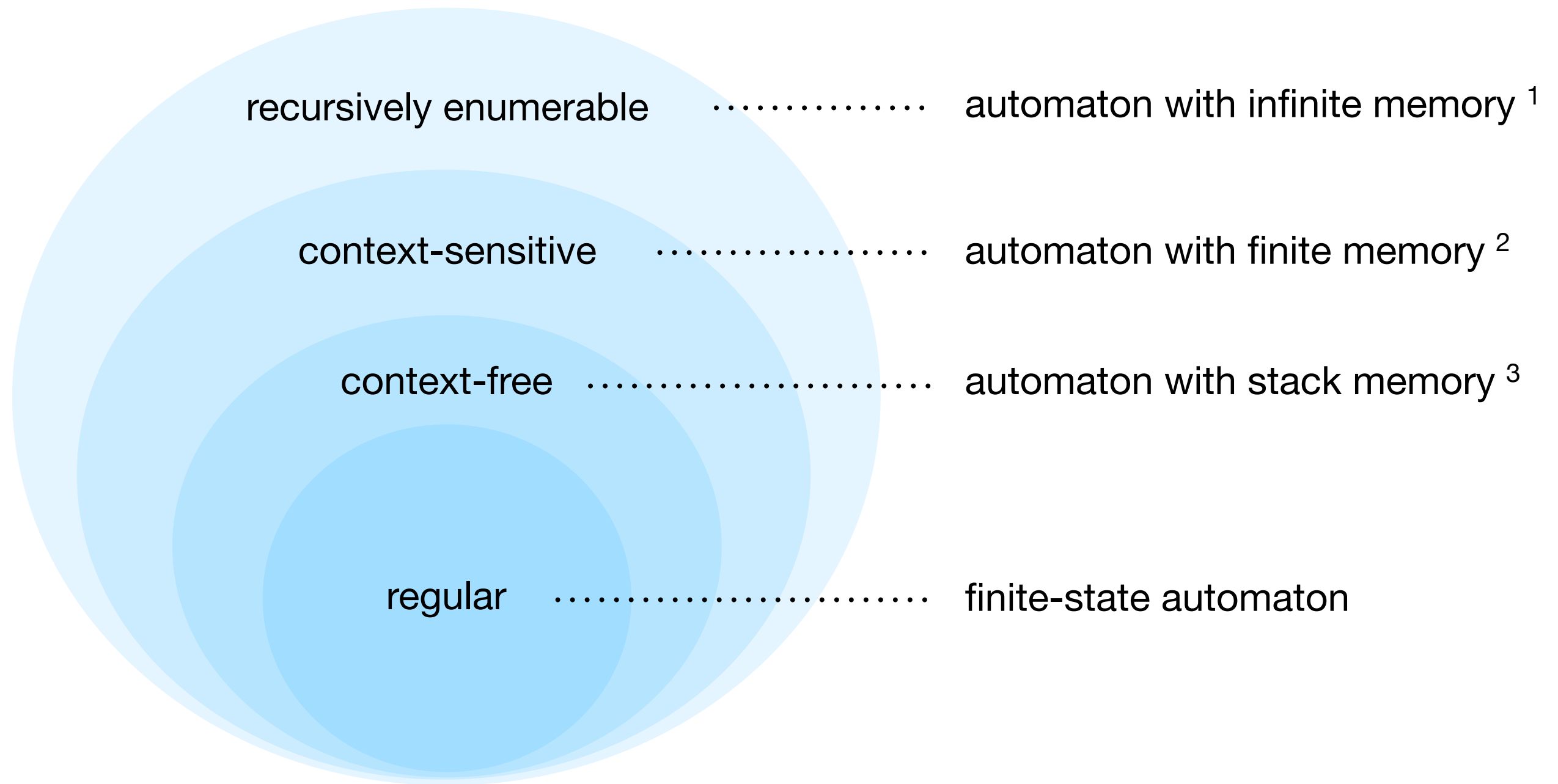


complex & rich



simple & limited

Chomsky hierarchy (1950's)



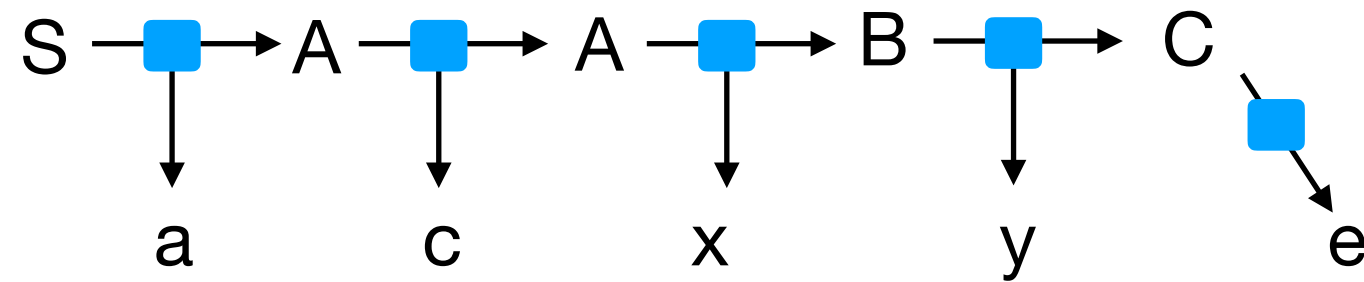
¹ Turing machine

² linear-bounded non-deterministic Turing machine

³ non-deterministic pushdown automaton

structure of derivations

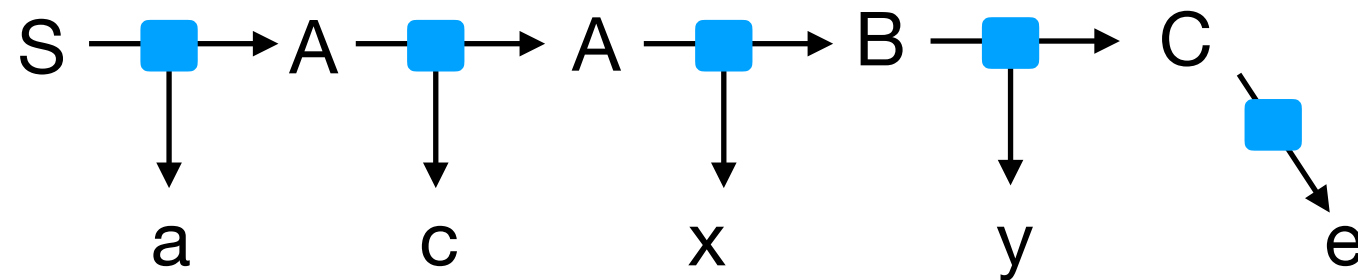
regular grammar:



- always linear
- used in computer science (e.g. search patterns)

structure of derivations

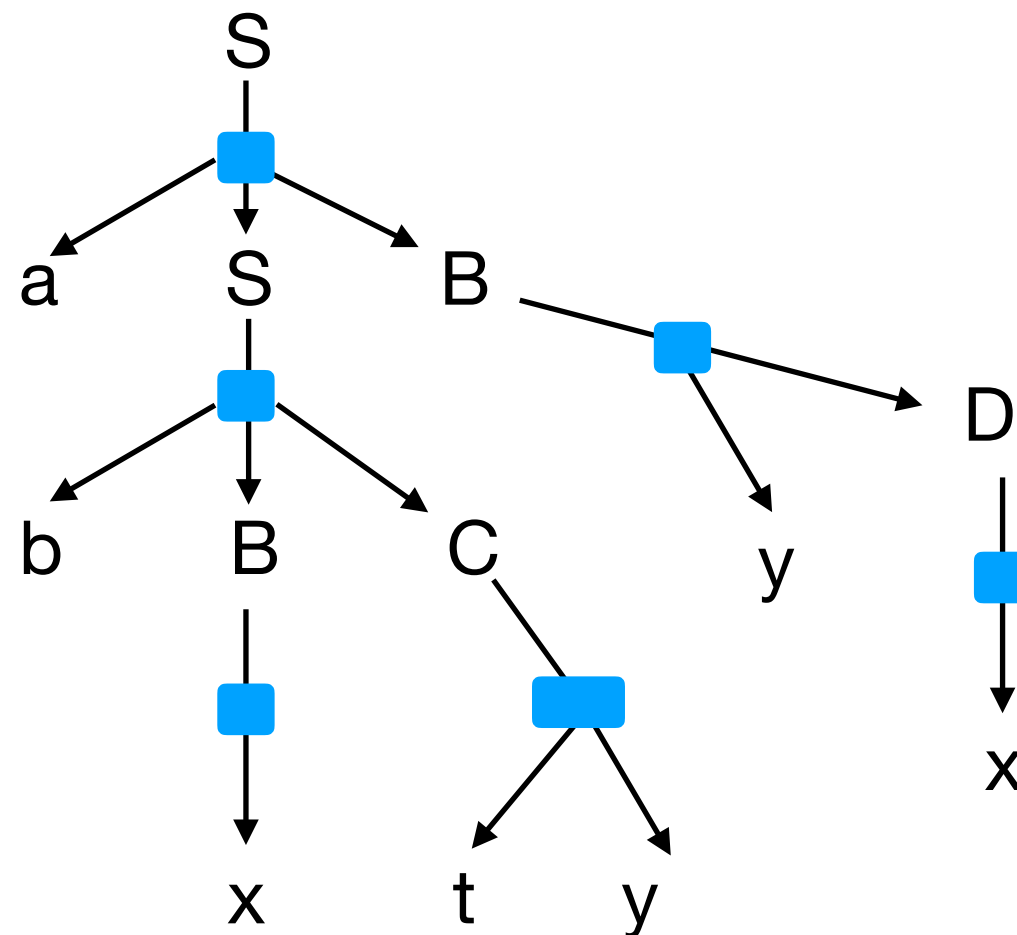
regular grammar:



- always linear
- used in computer science (e.g. search patterns)

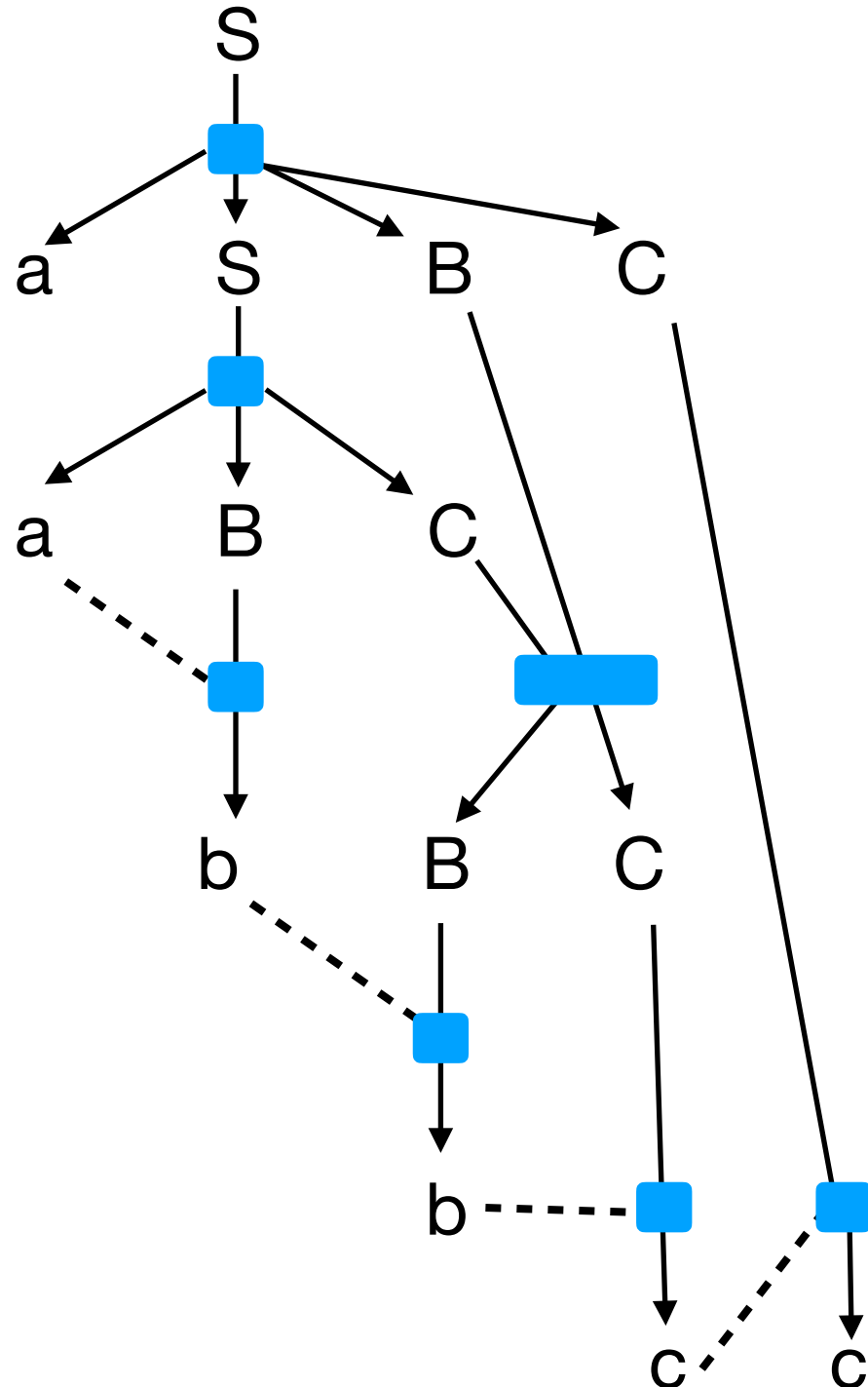
context-free grammar:

- always a tree
- used in linguistics for phrase structure (Chomsky 1956)
- central to computer science since Backus-Naur works ~1960



structure of derivations

context-sensitive grammar:



$S \Rightarrow aSBC \Rightarrow aaBCBC \Rightarrow aabCBC$
 $\Rightarrow aabBCC \Rightarrow aabbCC \Rightarrow aabbccC$
 $\Rightarrow aabbcc$

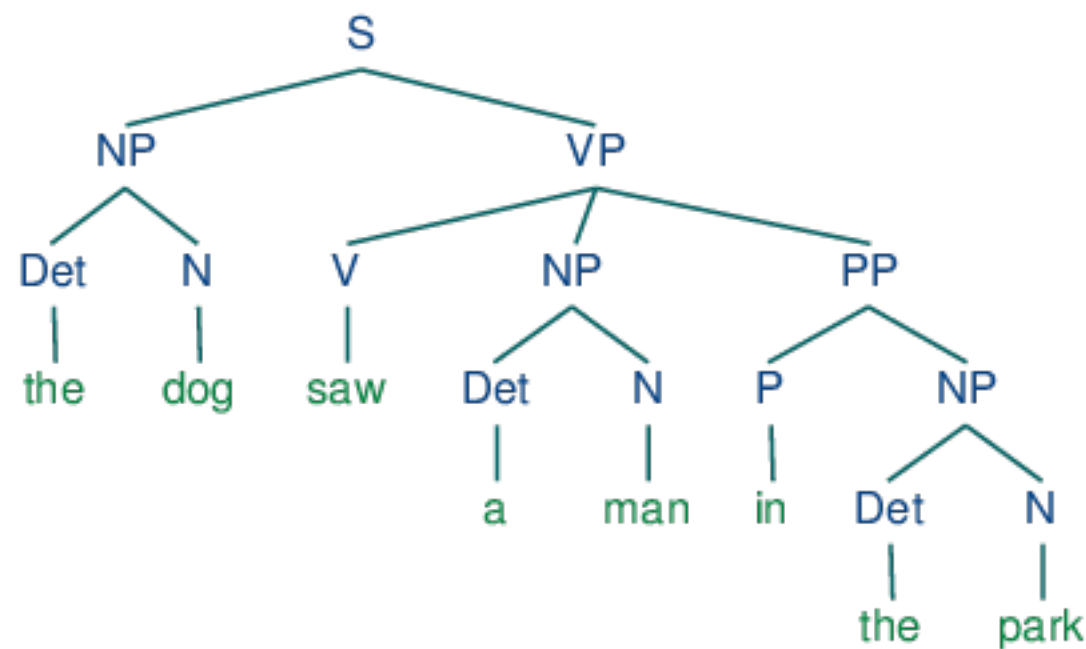
grammar:

$S \rightarrow aSBC$
 $S \rightarrow aBC$
 $CB \rightarrow BC$
 $aB \rightarrow ab$
 $bB \rightarrow bb$
 $bC \rightarrow bc$
 $cC \rightarrow cc$

what about natural languages?

- ~7000 existing languages
- only 2 have confirmed non-context-free features (Swiss-German, Bambara)

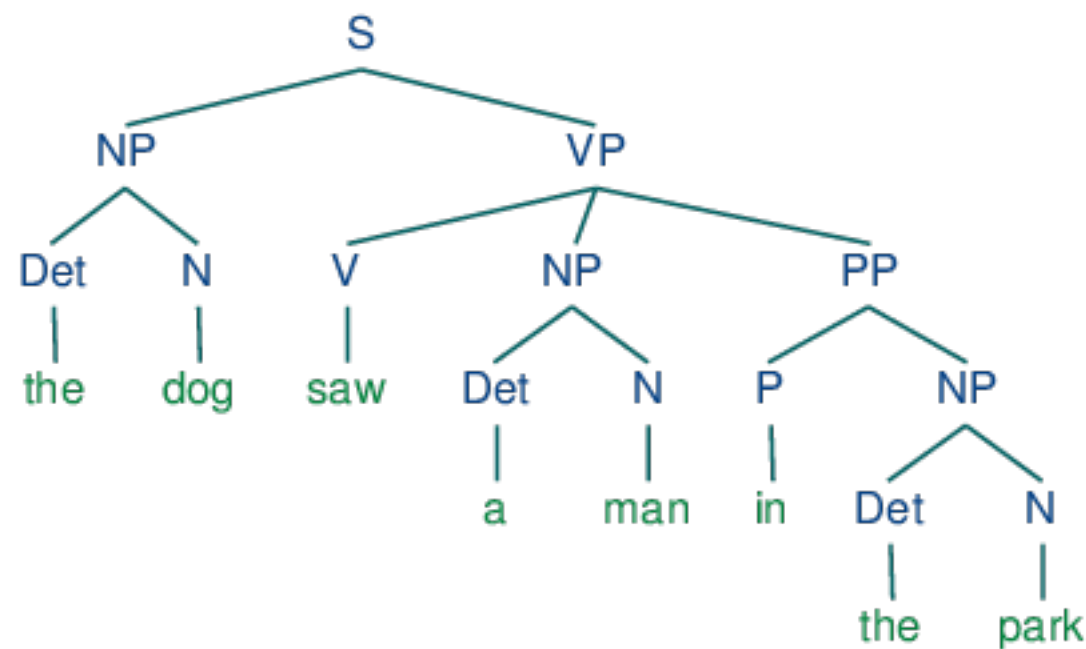
i.e. context-free languages define an *ensemble* for natural language syntax



what about natural languages?

- ~7000 existing languages
- only 2 have confirmed non-context-free features (Swiss-German, Bambara)

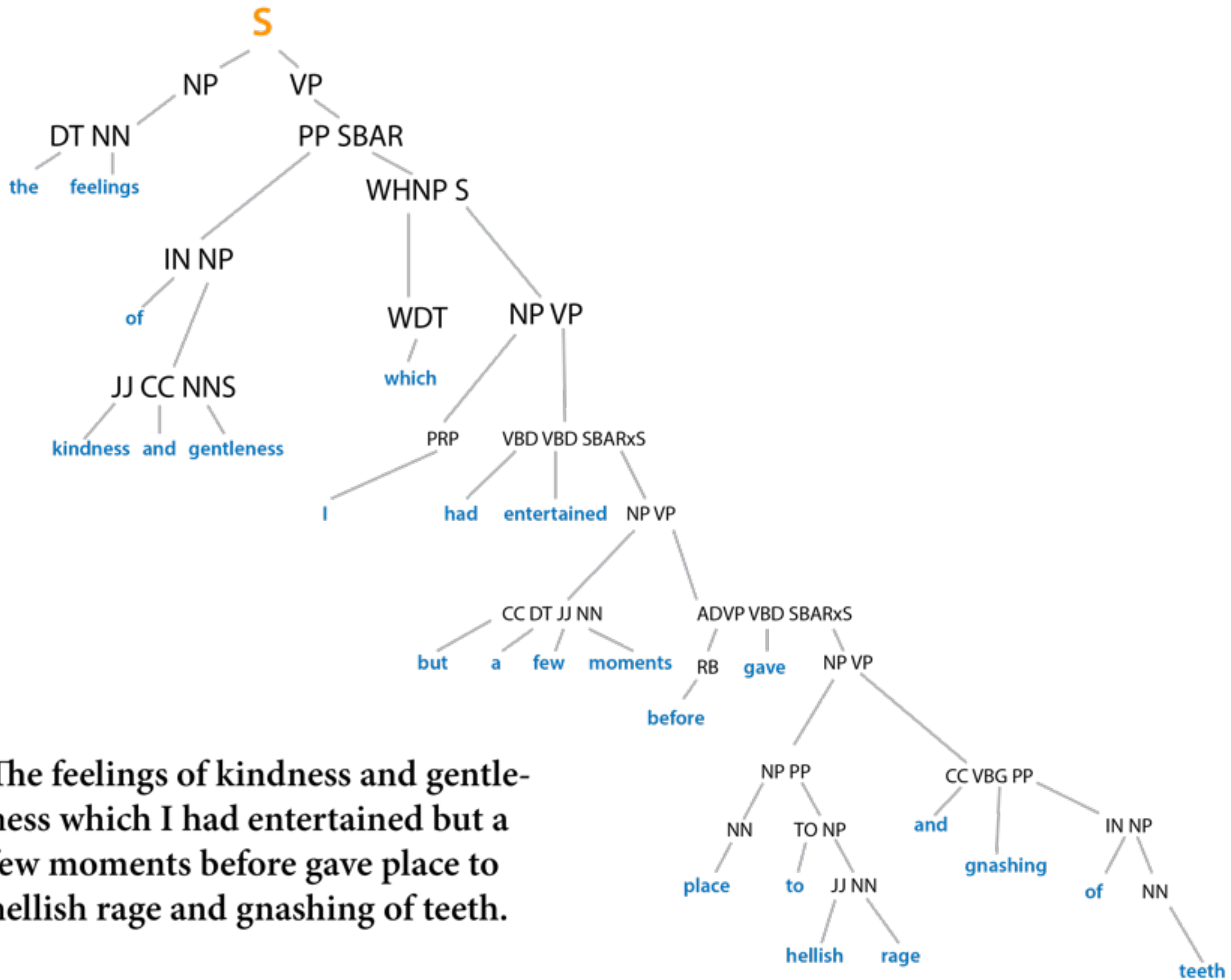
i.e. context-free languages define an *ensemble* for natural language syntax



meaning of the tree?

‘the park’ behaves like ‘park’

‘in the park’ behaves like ‘in—noun’



The feelings of kindness and gentleness which I had entertained but a few moments before gave place to hellish rage and gnashing of teeth.

2nd principle of language: **hierarchy**

The feelings of kindness and gentleness which I had entertained but a few moments before gave place to hellish rage and gnashing of teeth.

```

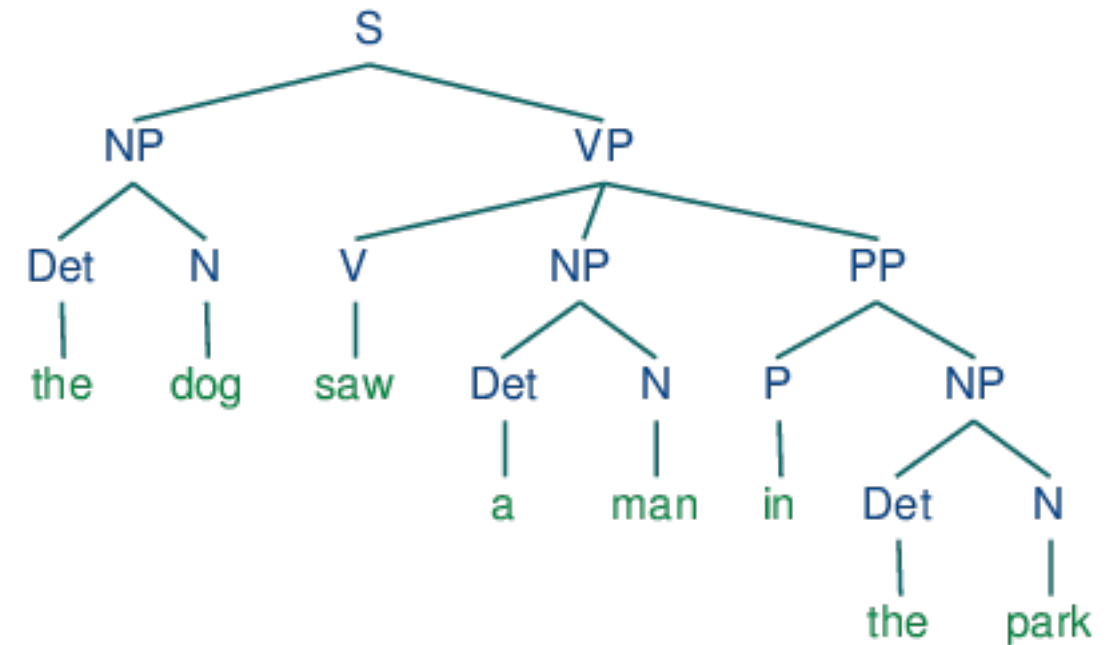
graph TD
    Root[ADVP VBD SBARxS] --- RB[RB]
    Root --- VBD[VBD]
    Root --- NP_VP[NP VP]
    RB --- Dogs[Dogs]
    VBD --- gave[gave]
    NP_VP --- NP1[NP]
    NP_VP --- VP[VP]
    NP1 --- NP_PP[NP PP]
    NP_PP --- NP2[NP]
    NP_PP --- PP1[PP]
    NP2 --- NN1[NN]
    NN1 --- place[place]
    PP1 --- TO[TO]
    TO --- to[to]
    PP1 --- NP3[NP]
    NP3 --- JJ[JJ]
    JJ --- hellish[hellish]
    NP3 --- NN2[NN]
    NN2 --- rage[rage]
    VP --- CC[CC]
    CC --- and[and]
    VP --- VBG_PP[VBG PP]
    VBG_PP --- VBG[VBG]
    VBG --- gnashing[gnashing]
    VBG_PP --- PP2[PP]
    PP2 --- IN[IN]
    IN --- of[of]
    PP2 --- NP4[NP]
    NP4 --- NN3[NN]
    NN3 --- teeth[teeth]
  
```

random language model

can we understand something about *typical* context-free grammars?

1. can assume binary tree¹

all rules either $A \rightarrow BC$ or $A \rightarrow b$



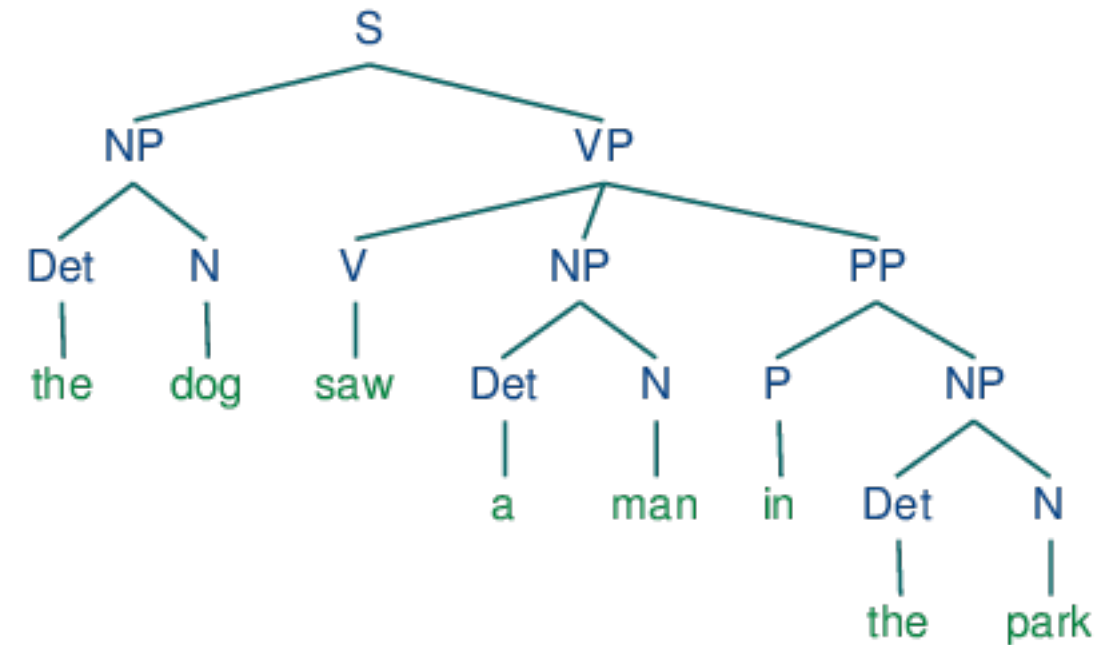
¹ binary tree = 'Chomsky normal form'

random language model

can we understand something about *typical* context-free grammars?

1. can assume binary tree¹

all rules either $A \rightarrow BC$ or $A \rightarrow b$



2. so far, rules have been yes/no. let rules \rightarrow conditional probabilities

then a grammar is defined by

$$M_{ABC} = \mathbb{P}(A \rightarrow BC \mid A \rightarrow \text{hidden}),$$

$$O_{Ab} = \mathbb{P}(A \rightarrow b \mid A \rightarrow \text{observable}),$$

¹ binary tree = 'Chomsky normal form'

random language model

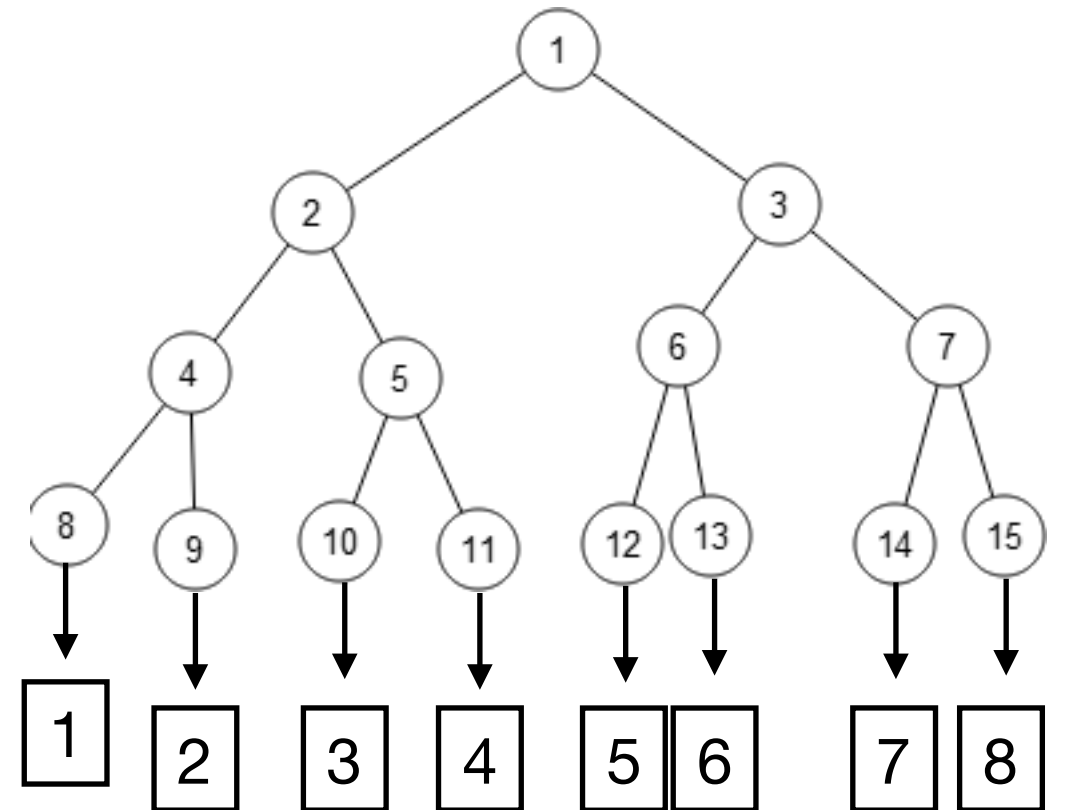
for simplicity, fix tree topology

$$M_{ABC} = \mathbb{P}(A \rightarrow BC \mid A \rightarrow \text{hidden}),$$

$$O_{Ab} = \mathbb{P}(A \rightarrow b \mid A \rightarrow \text{observable}),$$

$$M_{\sigma_i \sigma_j \sigma_k} = \mathbb{P}(\sigma_i \rightarrow \sigma_j \sigma_k \mid \sigma_i, \sigma_j, \sigma_k \in \chi_N),$$

$$O_{\sigma_i o_j} = \mathbb{P}(\sigma_i \rightarrow o_j \mid \sigma_i \in \chi_N, o_j \in \chi_T),$$



$$\mathbb{P}(\{\sigma_i, o_t\} \mid M, O, \mathcal{T}) = P_{\sigma_0} \prod_{\alpha \in \Omega} M_{\sigma_{\alpha_1} \sigma_{\alpha_2} \sigma_{\alpha_3}} \prod_{\alpha \in \partial \Omega} O_{\sigma_{\alpha_1} o_{\alpha_2}},$$

note: M,O are probabilities for a fixed grammar, then we have an ensemble of grammars

random language model

what is the physics?

$$M_{\sigma_i \sigma_j \sigma_k} = \mathbb{P}(\sigma_i \rightarrow \sigma_j \sigma_k | \sigma_i, \sigma_j, \sigma_k \in \chi_N),$$
$$O_{\sigma_i o_j} = \mathbb{P}(\sigma_i \rightarrow o_j | \sigma_i \in \chi_N, o_j \in \chi_T),$$

$$Z = \int DM \int DO \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{\log \mathbb{P}}$$

Z contains all the context-free grammars
& all grammatical sentences in the
universe!

random language model

what is the physics?

$$M_{\sigma_i \sigma_j \sigma_k} = \mathbb{P}(\sigma_i \rightarrow \sigma_j \sigma_k | \sigma_i, \sigma_j, \sigma_k \in \chi_N),$$
$$O_{\sigma_i o_j} = \mathbb{P}(\sigma_i \rightarrow o_j | \sigma_i \in \chi_N, o_j \in \chi_T),$$

$$Z = \int DM \int DO \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{\log \mathbb{P}}$$

impose normalization of probabilities
& # of nonzero rules

Z contains all the context-free grammars
& all grammatical sentences in the
universe!

random language model

what is the physics?

$$M_{\sigma_i \sigma_j \sigma_k} = \mathbb{P}(\sigma_i \rightarrow \sigma_j \sigma_k | \sigma_i, \sigma_j, \sigma_k \in \chi_N),$$

$$O_{\sigma_i o_j} = \mathbb{P}(\sigma_i \rightarrow o_j | \sigma_i \in \chi_N, o_j \in \chi_T),$$

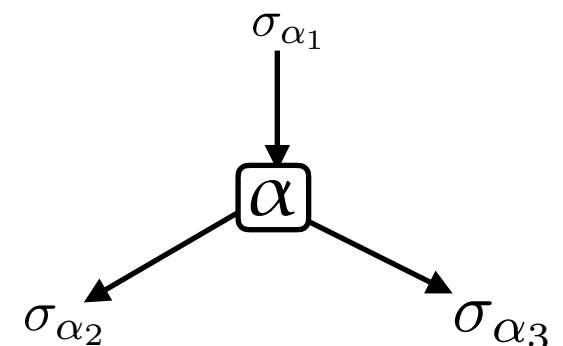
$$Z = \int DM \int DO \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{\log \mathbb{P}}$$

impose normalization of probabilities
& # of nonzero rules

Z contains all the context-free grammars
& all grammatical sentences in the
universe!

$$\log \mathbb{P}(\{\sigma_i, o_t\} | M, O, \mathcal{T}) = \log P_{\sigma_0} + \sum_{\alpha \in \Omega} \log M_{\sigma_{\alpha_1} \sigma_{\alpha_2} \sigma_{\alpha_3}} + \sum_{\alpha \in \partial \Omega} \log O_{\sigma_{\alpha_1} o_{\alpha_2}}$$

looks a bit like a spin model ... except



random language model

remarkably we can count all the context-free grammars in the universe

$$Z = \int DM \int DO \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{\log \mathbb{P}}$$

what is the miracle? discrete Fourier transform



$$Z = Z_0 \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{-H}$$

random language model

remarkably we can count all the context-free grammars in the universe

$$Z = \int DM \int DO \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{\log \mathbb{P}}$$

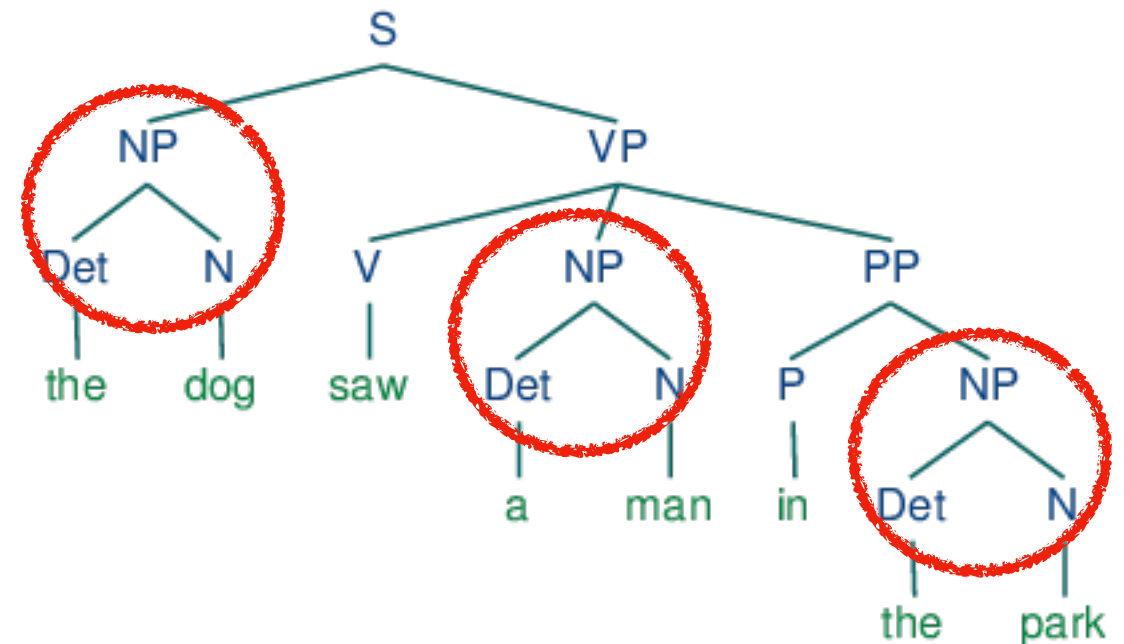
what is the miracle? discrete Fourier transform



$$H = -k \sum_{\alpha, \beta \in \Omega} (N^2 \delta_{\sigma|\alpha, \sigma|\beta} - \delta_{\sigma_{\alpha_1}, \sigma_{\beta_1}}) + \dots$$

~ Potts model

$$Z = Z_0 \sum_{\mathcal{T}} \sum_{\{\sigma\}} \sum_{\{o\}} e^{-H}$$



the SETI problem



ଫର୍ଷ୍ଟଫଗଜପଞ୍ଜିଫର୍ଷ୍ଟଫକ୍ଷପଦଗଓଡ଼ବର୍ଷଓଠଦଡ଼କଦବଗଡ଼ବରଗଘପଡ଼କଶଡ଼ଜପଶଞଠଦଧରବଓ
ଭଦକ୍ଷଫନଧପଡ଼କଠଠବବଫଭଞଞଓପଡ଼କ୍ଷଧଧଡ଼ଜରଫଠଦଡ଼କକଓବଶଜଜଗଠରଜରଜଦଘପଟଟଘ
ଗଜଫଘଟଟଶଗଶକଡ଼ଡ଼ଭର୍ଜରଠଡ଼ଓଟଟଧଭଶଡ଼ଶଟଘଓକ୍ଷଓଦନଟଫଡ଼ଟକଓଘକଡ଼ଡ଼କ୍ଷଓଗଶଜକ
ଗଶଧଧଡ଼ବଡ଼କ୍ଷଡ଼କ୍ଷଡ଼କଟଭଓଧଶଠଶଘଡ଼ଜଫଧଭନକବପଶଶଫର୍ଷ୍ଟଜକଡ଼ଘଜଜବଧବଗବଘଘ
ଡ଼ଘଶଦଡ଼ଗଧଟବବଦପଦଫଓ

can we learn its language?

the SETI problem

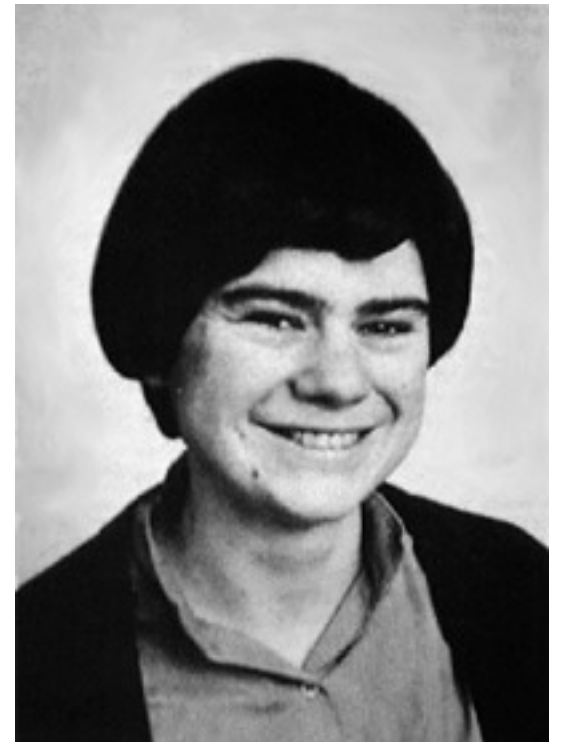
[illegible]

can we learn its language?

assume it is generated by a CFG

count number of grammars for which text is grammatical ¹

'Gardner' volume

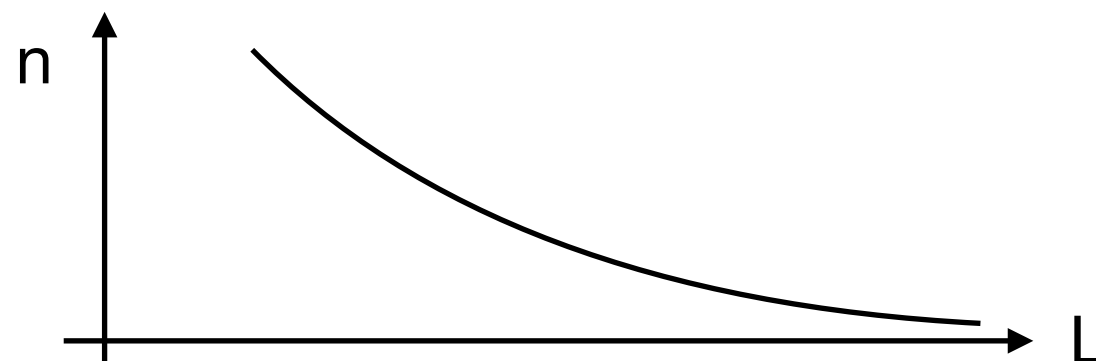
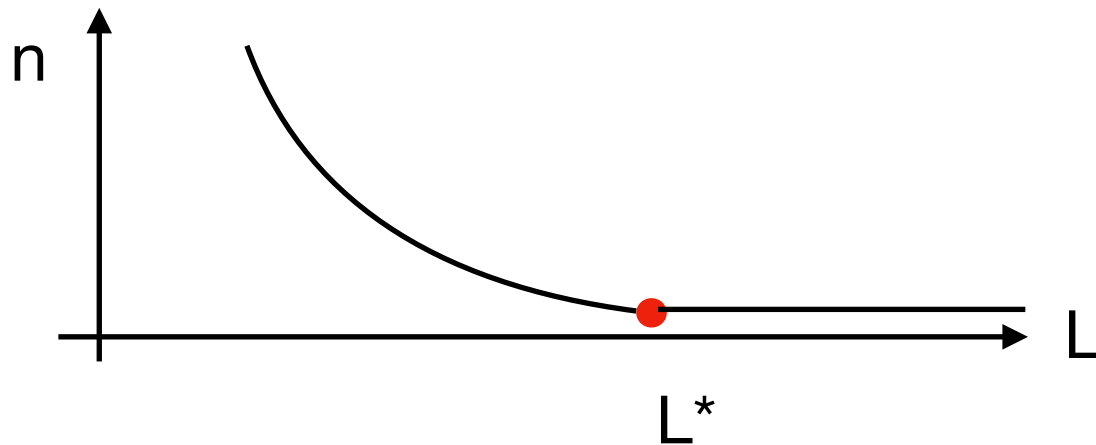


¹ more precisely, below some threshold K in probability

the SETI problem

[illegible]

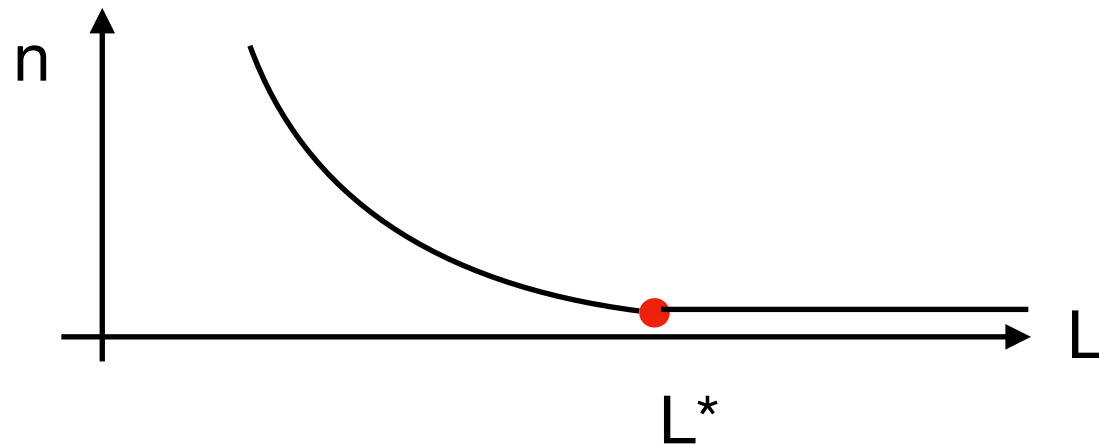
what do we expect?



the SETI problem

[illegible]

what do we expect?



I am working on the full solution..

in a simple (wrong) approximation
it is equivalent to Gardner's result
for the perceptron

$$L = \frac{\ell}{N_{DOF}}$$

perspectives

ambiguity: For a typical sentence, how many grammatical parses are there?

If $n = 1$, sentence is unambiguous

If $n > 1$, sentence is ambiguous

If $n = 0$, sentence is ungrammatical

Natural languages are typically ambiguous, e.g.

“Two cars were reported stolen by the Groverton police yesterday”¹

¹ from S Pinker, The Language Instinct

perspectives

ambiguity: For a typical sentence, how many grammatical parses are there?

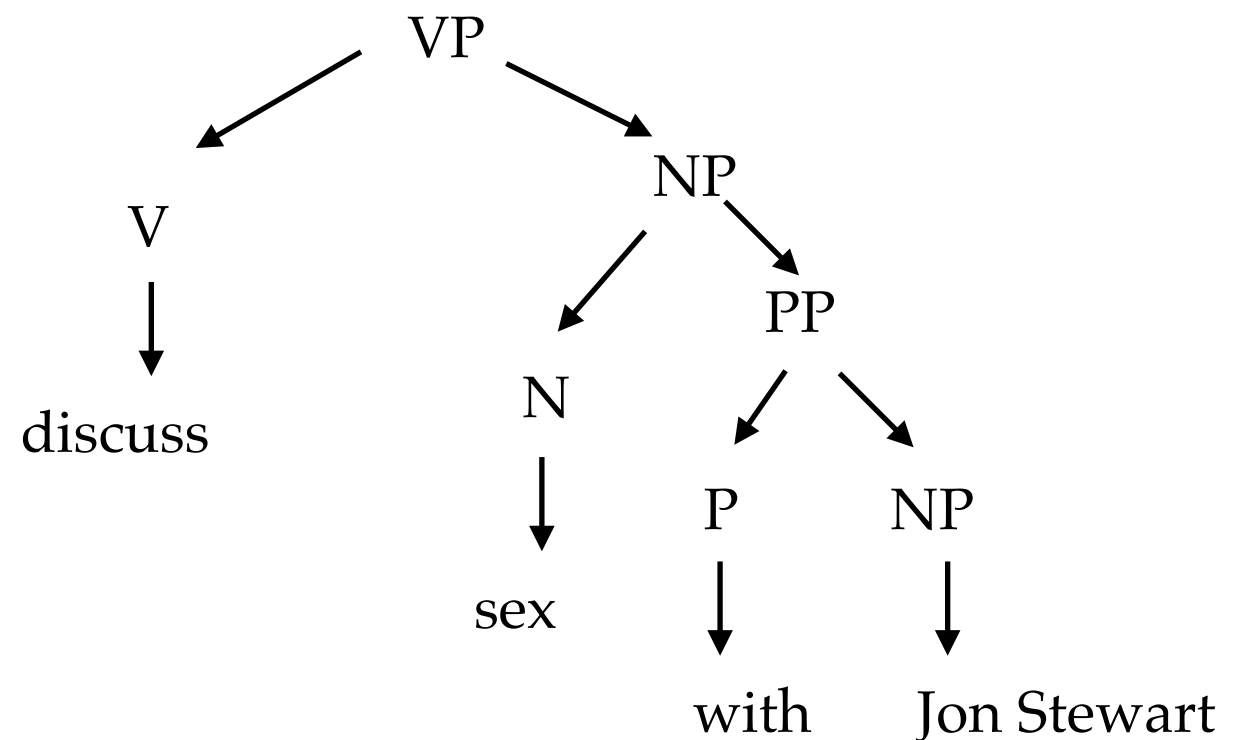
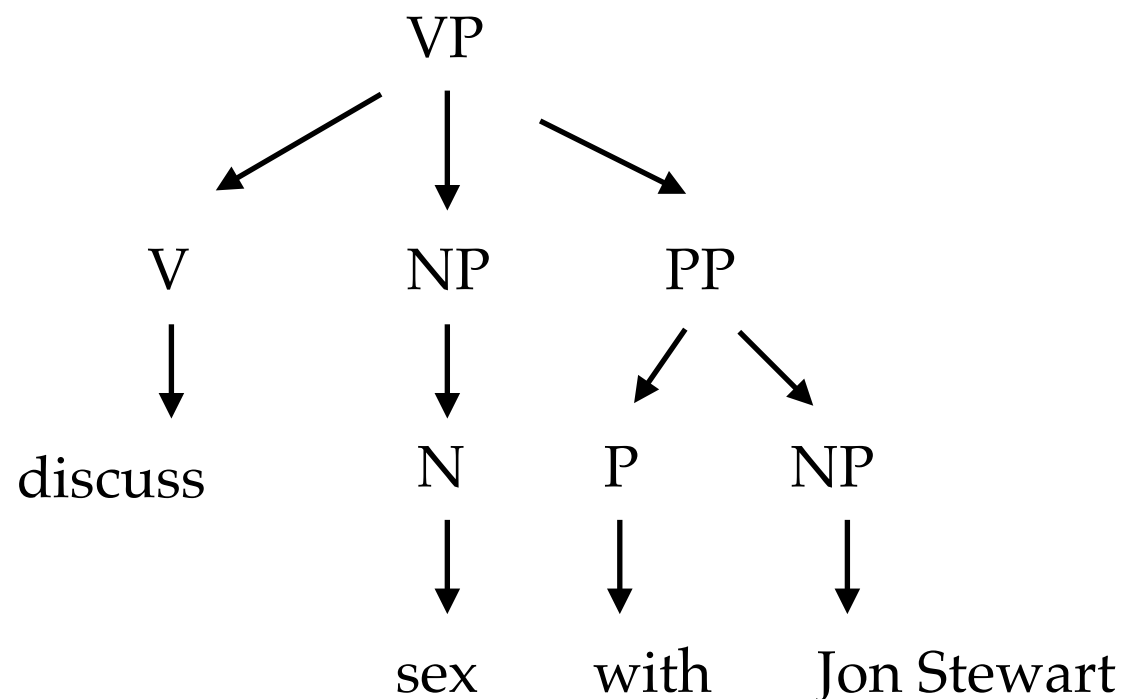
If $n = 1$, sentence is unambiguous

If $n > 1$, sentence is ambiguous

If $n = 0$, sentence is ungrammatical

Natural languages are typically ambiguous, e.g.

“Two cars were reported stolen by the Groverton police yesterday”¹



¹ from S Pinker, The Language Instinct

perspectives

ambiguity: For a typical sentence, how many grammatical parses are there?

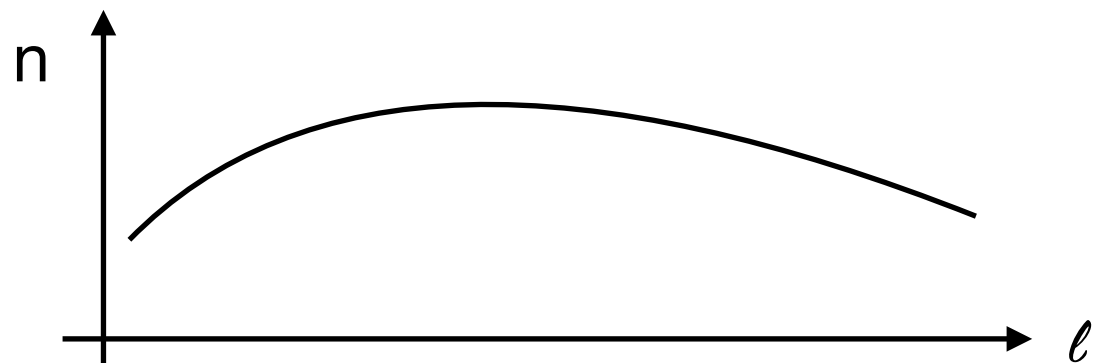
If $n = 1$, sentence is unambiguous

If $n > 1$, sentence is ambiguous

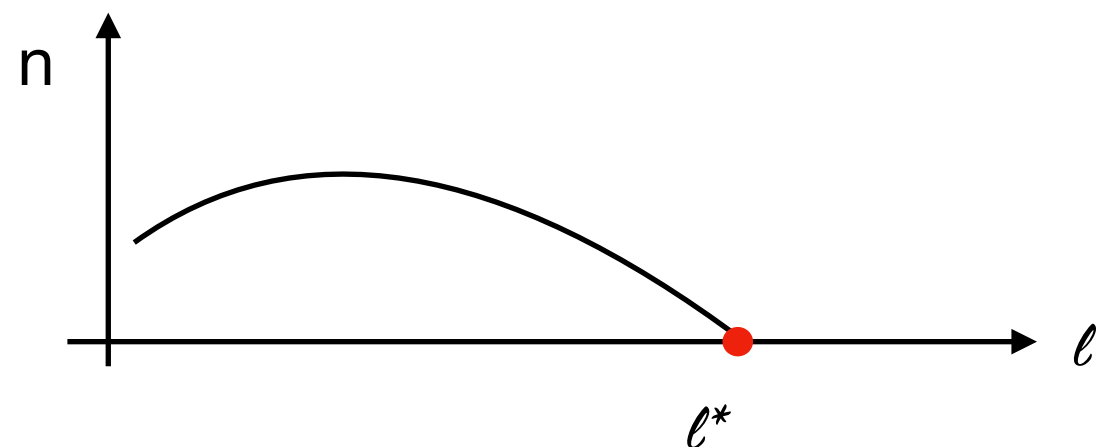
If $n = 0$, sentence is ungrammatical

Natural languages are typically ambiguous, e.g.

“Two cars were reported stolen by the Groverton police yesterday”¹



which situation?



¹ from S Pinker, The Language Instinct

perspectives

phase diagram: What is the phase diagram of languages?

Are human languages atypical?

neural networks and learning:

What is the optimal architecture to learn highly compositional functions?

For natural language processing, how best to incorporate syntax into neural network approaches?

Can tools of disordered physics (e.g Thouless-Anderson-Palmer equations) help to learn languages?

perspectives

semantics: syntax isn't everything..

e.g. who is 'he' in this dialogue: ¹

Alice: I'm leaving you.

Bob: Who is he?!

Is there a physical approach to semantics?

c.f. dependency grammars, Montague grammars, ...

¹ from S Pinker, The Language Instinct

conclusions



- successful machine learning architectures are *compositional* and *hierarchical*
- natural languages are also compositional and hierarchical
- context-free grammars define a simple model for these properties
- ensemble of grammars = random language model
- the statistical mechanical problem is not trivial, but not intractable

Mathematical linguistics has been around for 60 years.
It's time for physical linguistics!

Thanks to my colleagues at ENS and elsewhere in Paris:

Remi Monasson, Jorge Kurchan, Francesco Zamponi,
Guilhem Semerjian, Pierfrancesco Urbani