# General Principles of Human and Machine Learning

Lecture 7: Compression and resource constraints

David G. Nagy

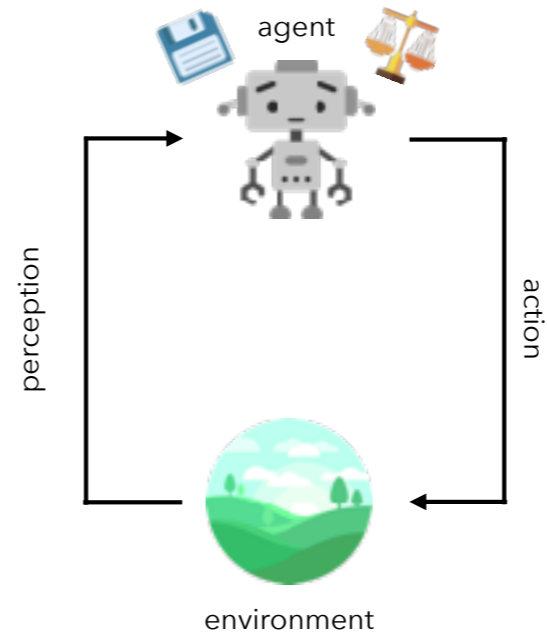# admin

- quiz
  - maximum score for quiz #2 is reduced to 16 from 18
  - from now on, pop quizzes may contain content from the lecture in the same week



- halfway feedback form on course, please submit at
  - https://forms.gle/BWBHobeVZniJKuKdA

**today**
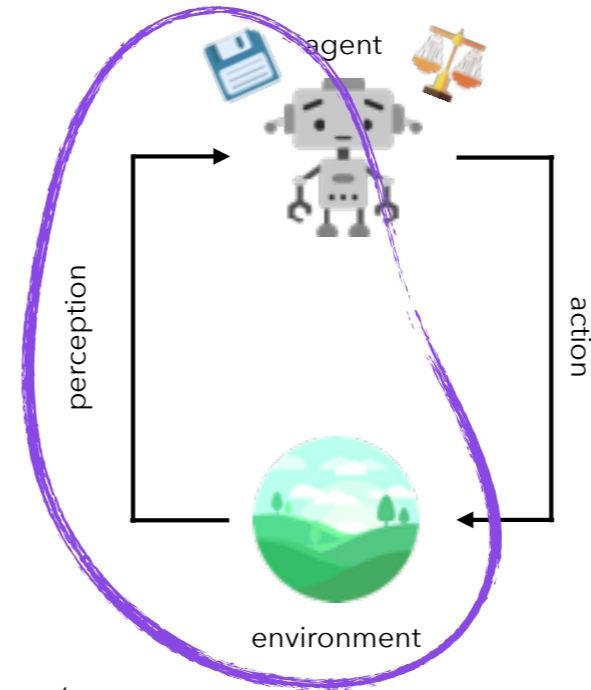
agent

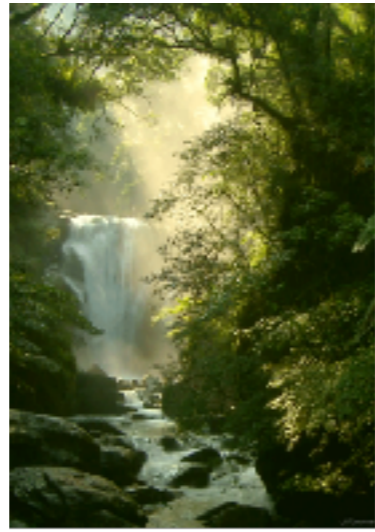perception

action

environment
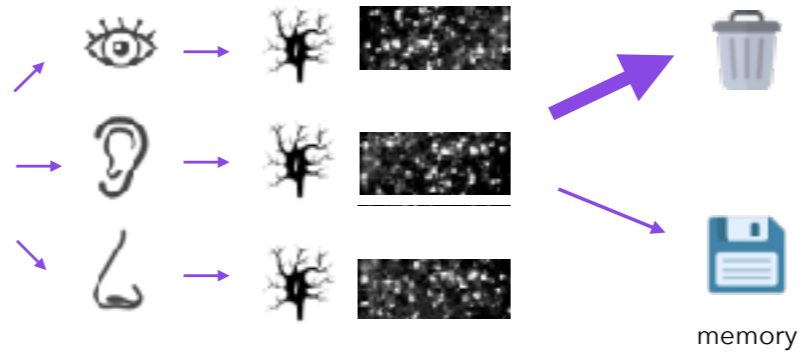
3

# today

- focus on constraints on **memory**

  - lossless and lossy compression

  - generative compression

    - perception as bayesian inference

    - human memory distortions

agent

perception

action

environment

4

# computational problem of memory



environment

(Nickerson & Adams, 1979)
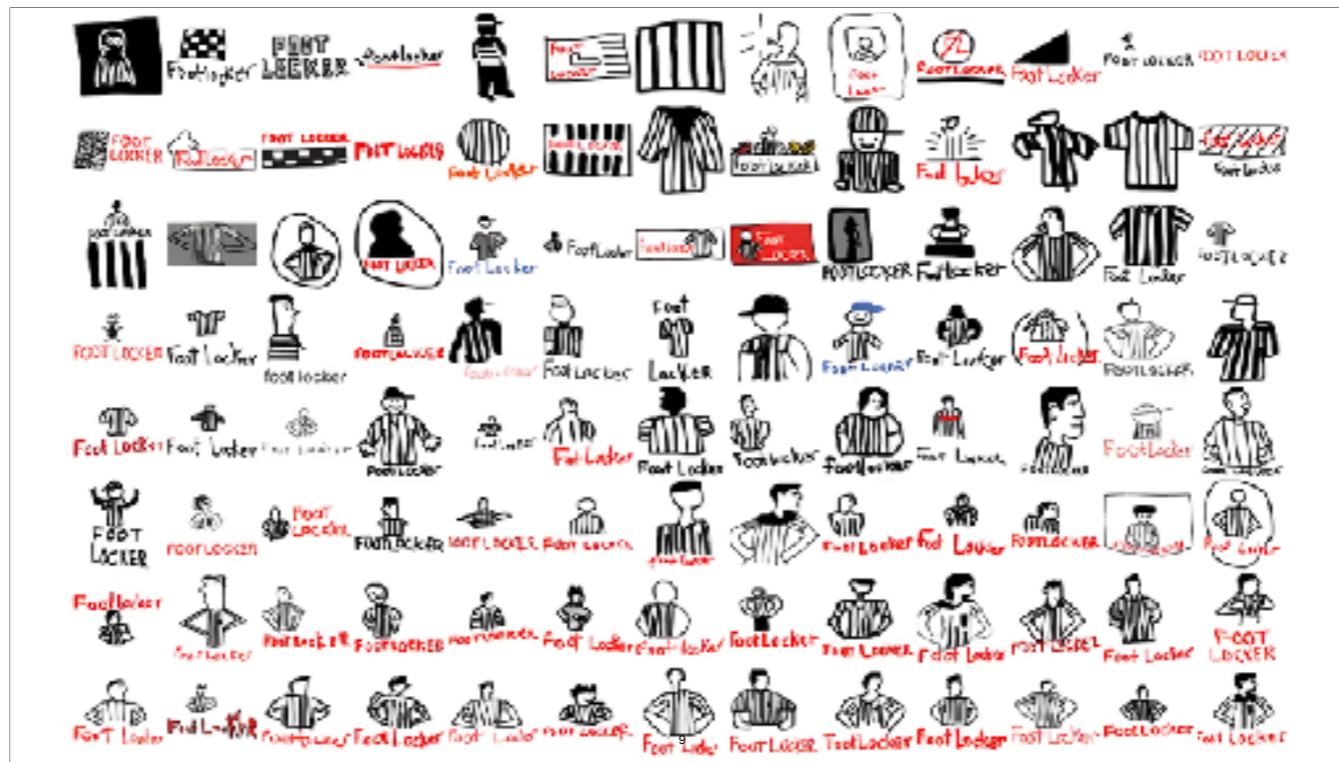
Also happens for logos of well known companies. Maybe people just aren't good at drawing?

near perfect
drawing ███ 8%

forgot referee ████████████████ 43%

referee facing the
wrong way ███████████████ 40%

added a hat ███████ 18%

shoe instead of
referee █████ 14%

# recognition



$P(A)$

$P(B)$

*A*

*B*

Prasad & Bainbridge, 2023

- is the lesson from this simply that human memory is poor?

- memory resources are certainly bounded

- but it is possible to do badly, do well or even optimally **in relation** to available resources

# compression

10101001101001101000010001000100101
11110100100110100100001011000011
10101101001101100000110111110110
10010101011101000110001000011000
10000000101001011001110010001000
01001100111000101100100101101010
10110110001000110001011000100001
11101101110000010110110011100100
11000100000011100111100101001001
01001010011111111111100110100111
01111110110111111110111000100111
10110110111011100010100101010010
10001000010100000000010101100000
00000110101001011010011001011000
00000000011001010101010100110011
10001100011000110100111110001001
100100110001101111100000111101001

<span>c</span> → 110101101110

13

theory of lossless compression

## compression

I walked my four legged animal that barks on the day before today after the huge glowing ball of fire left the sky

**compression**

I walked my four legged animal that barks on the day before today after the huge glowing ball of fire left the sky

## compression

I walked my dog on the day before today after the huge glowing ball of fire left the sky

**compression**

I walked my dog yesterday after the huge glowing ball of fire left the sky

## compression

I walked my dog yesterday after sunset

possible
inputs

**c**

00000000000 →

00000000001 → how should
we choose
00000000011 → the codes?

⋮

00001001101 → frequent inputs
↓
short codes

⋮

11111111111 →

19

frequency

possible
inputs

**c**
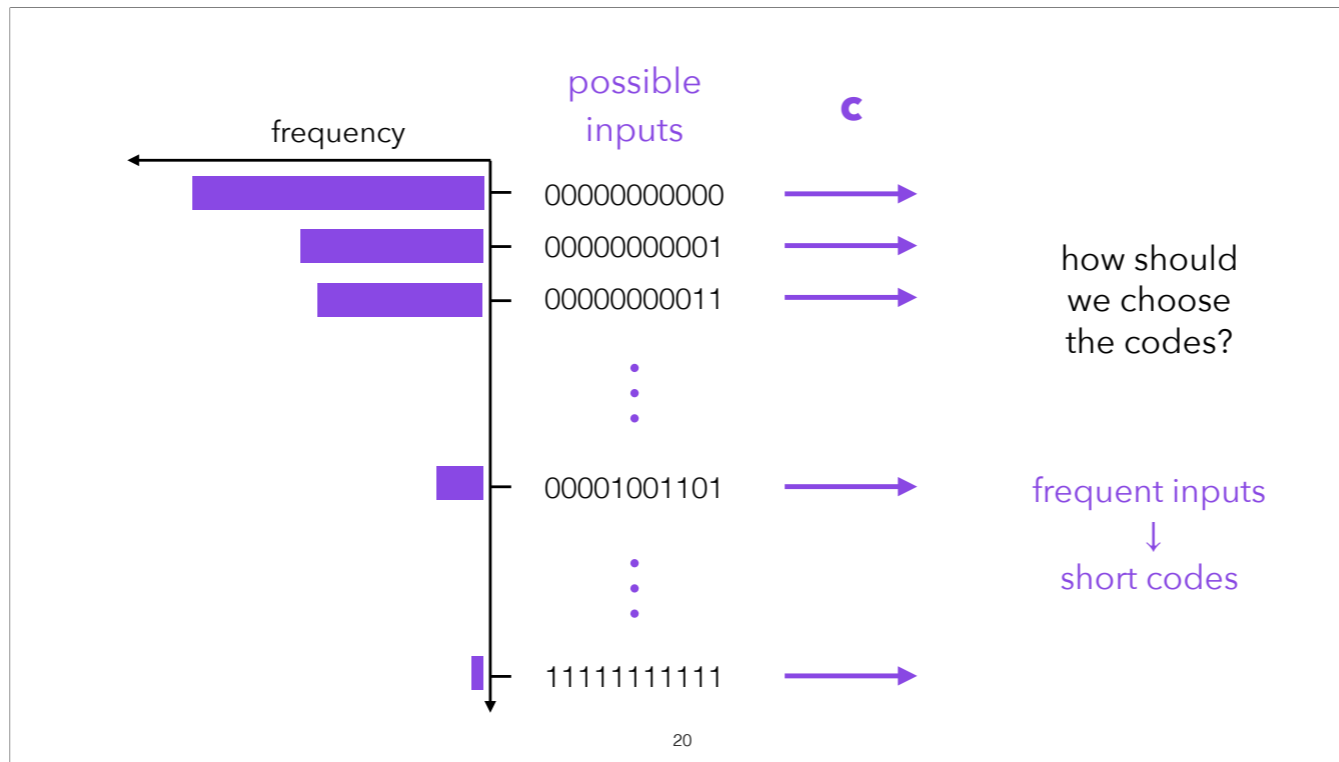
00000000000

00000000001

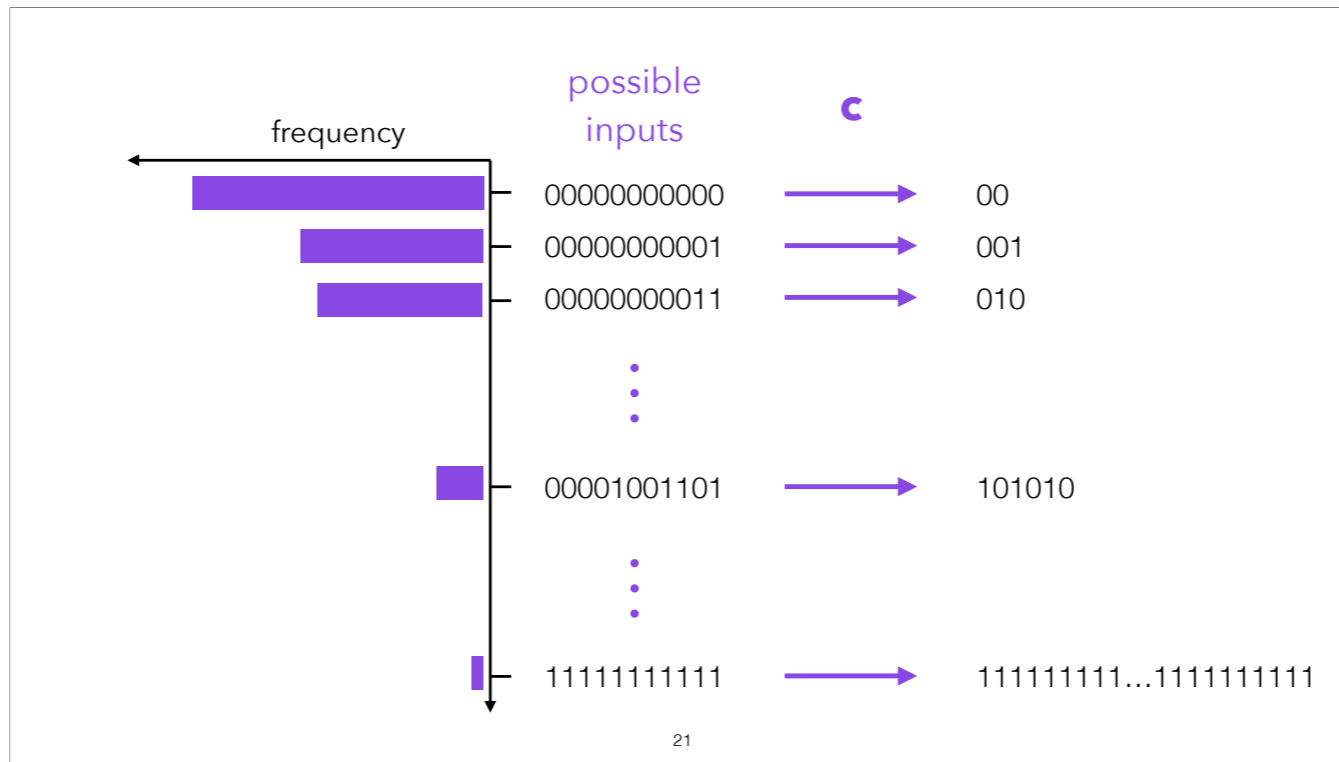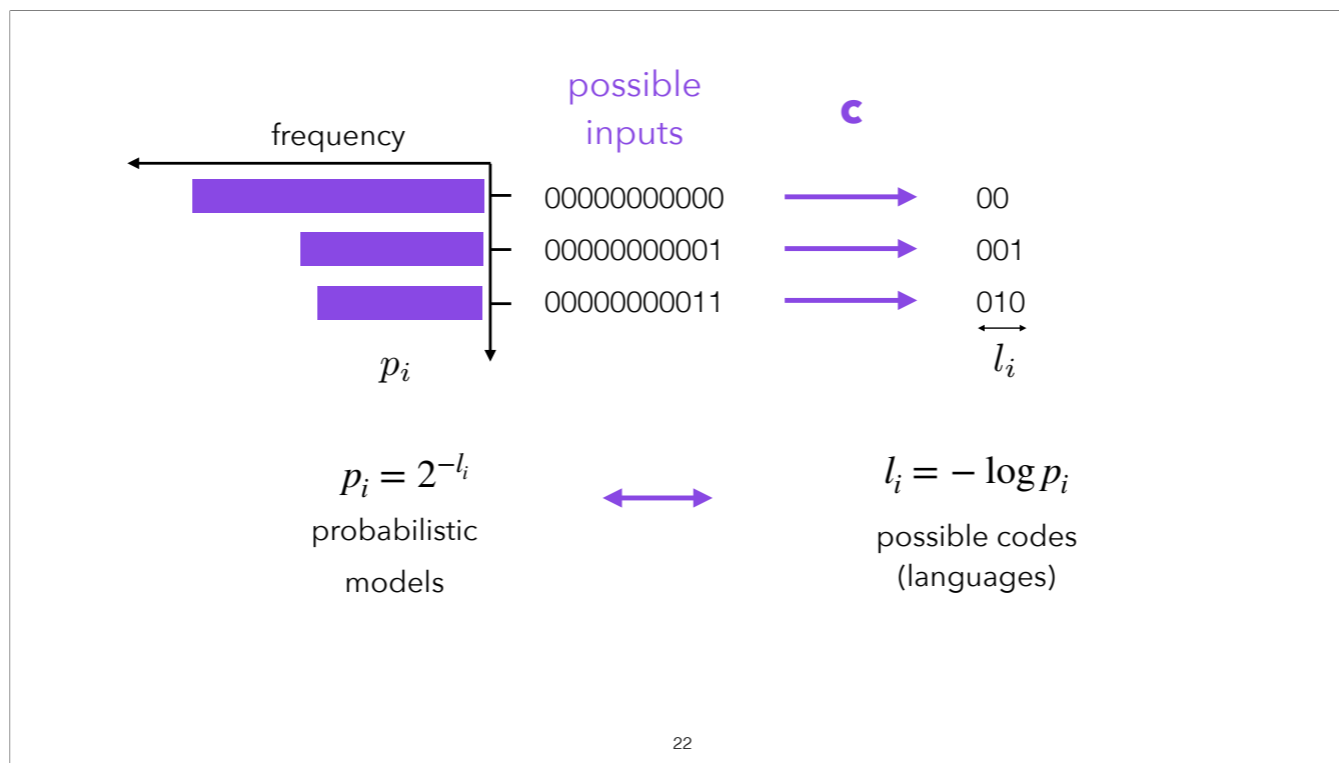00000000011

⋮

00001001101

⋮

11111111111

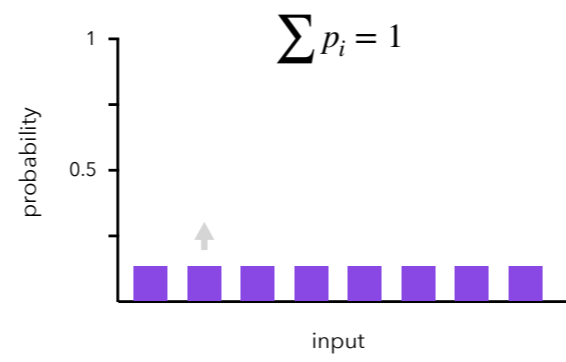how should
we choose
the codes?

frequent inputs
↓
short codes

note that using short codewords for frequent inputs means that the codewords for some inputs will have to be longer than the original input, so if our frequency estimates are wrong, the encoding might turn out to require more memory resources than just storing the original input directly
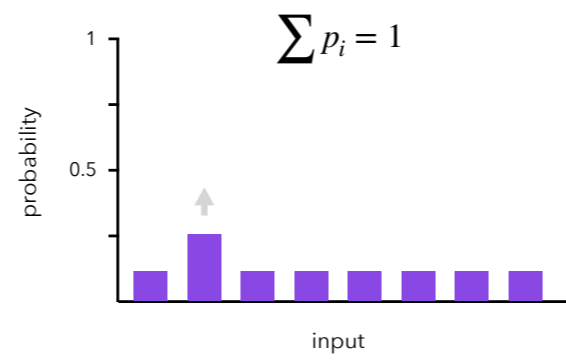
p_i is the probability or relative frequency of an input string

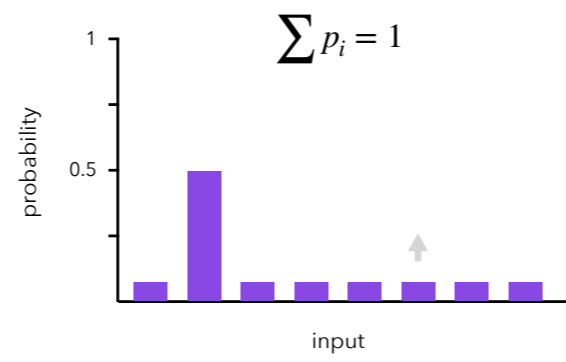l_i is the length of the code word that is used for the i-th input
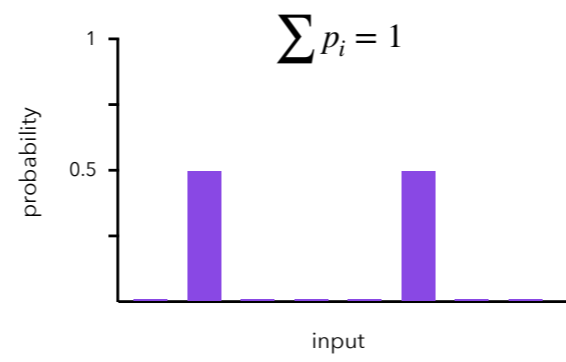
# probability budget



$$\sum p_i = 1$$

probability axis labeled with 1, 0.5; horizontal axis labeled "input"

$$\sum p_i = 1$$

$$\sum p_i = 1$$

$$\sum p_i = 1$$

# codeword budget

how many codewords of length $l_i$?

1

0

how many codewords of length $l_i$?

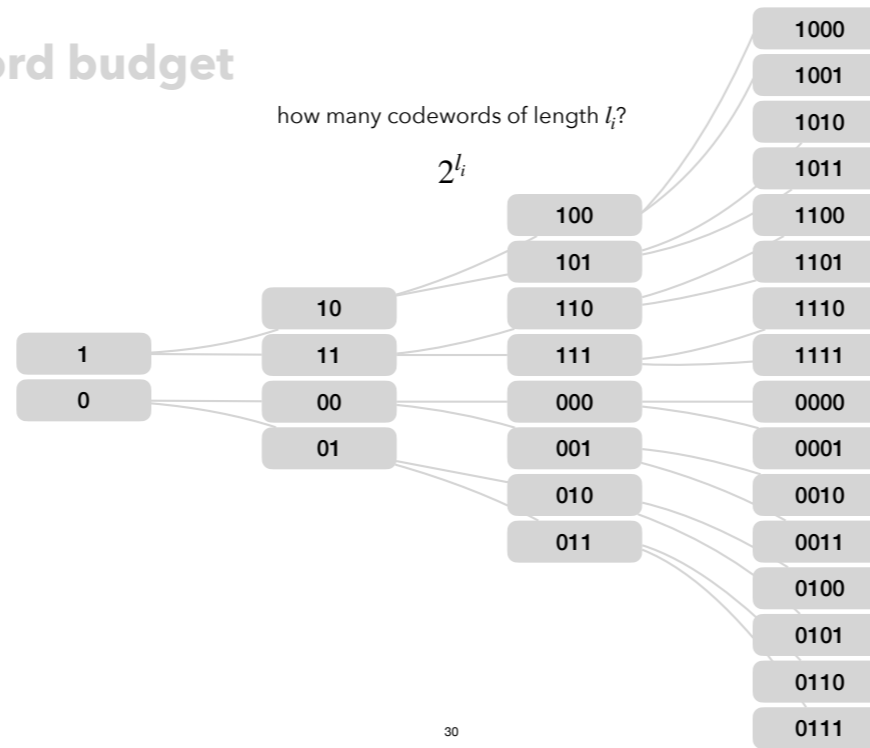| | |
|---|---|
| | 10 |
| 1 | 11 |
| 0 | 00 |
| | 01 |

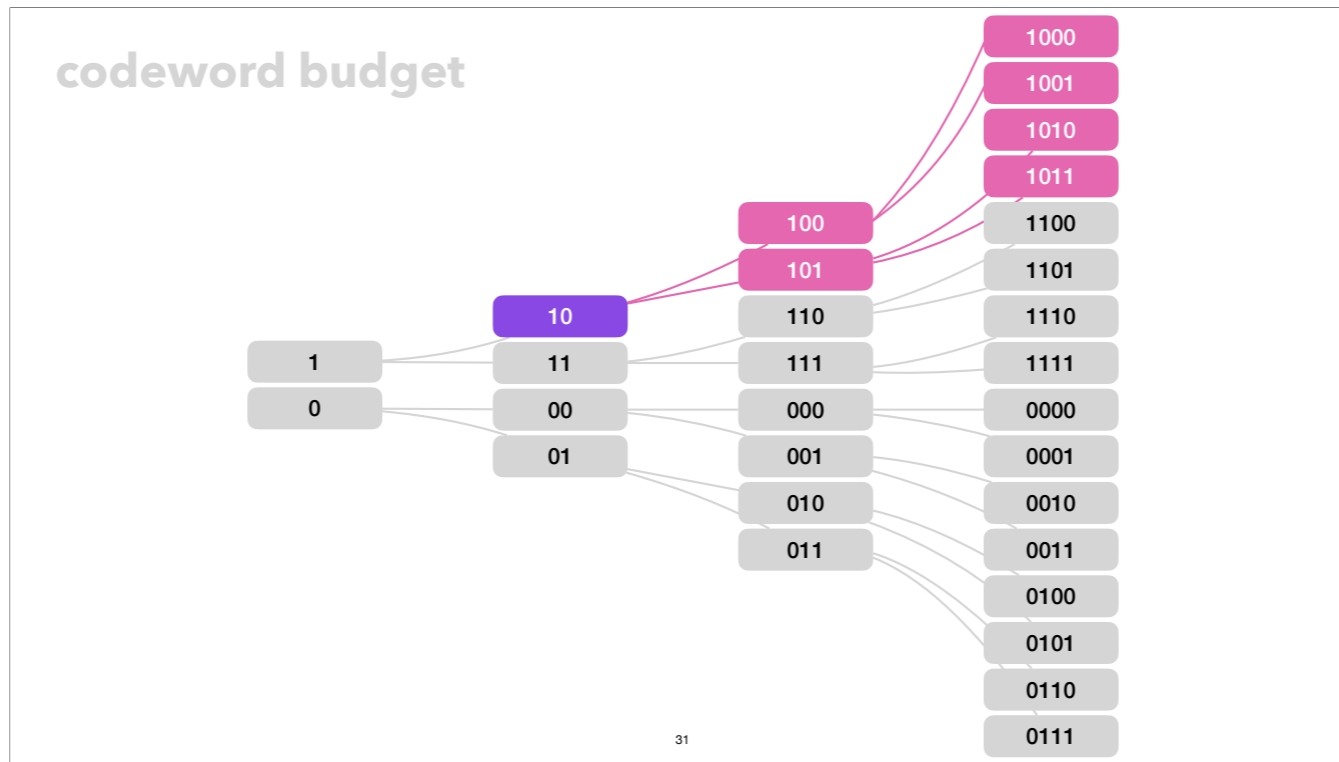# codeword budget

how many codewords of length $l_i$?

# codeword budget

how many codewords of length $l_i$?

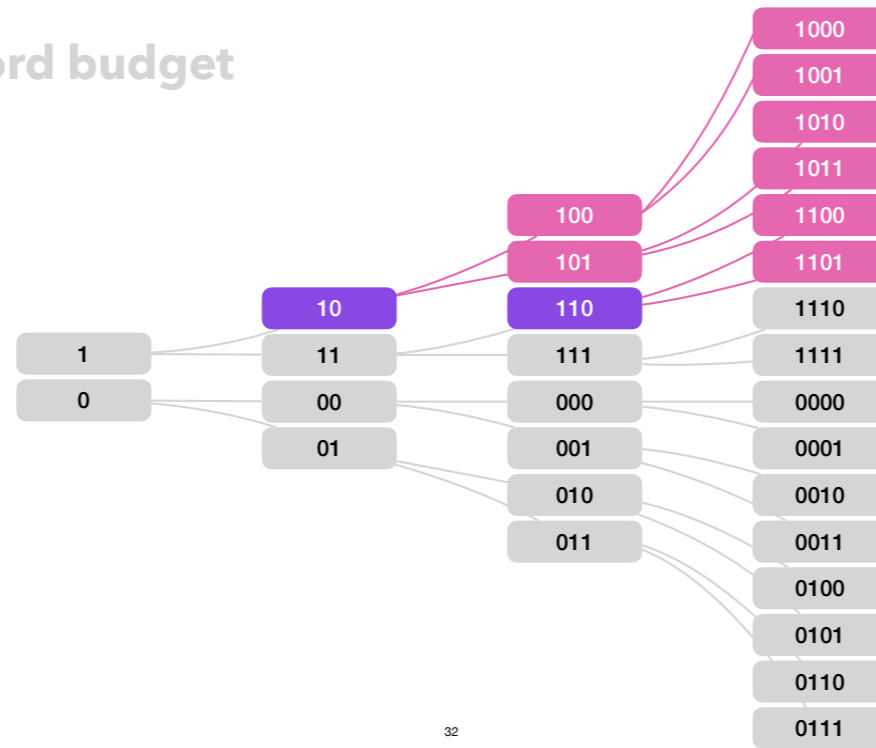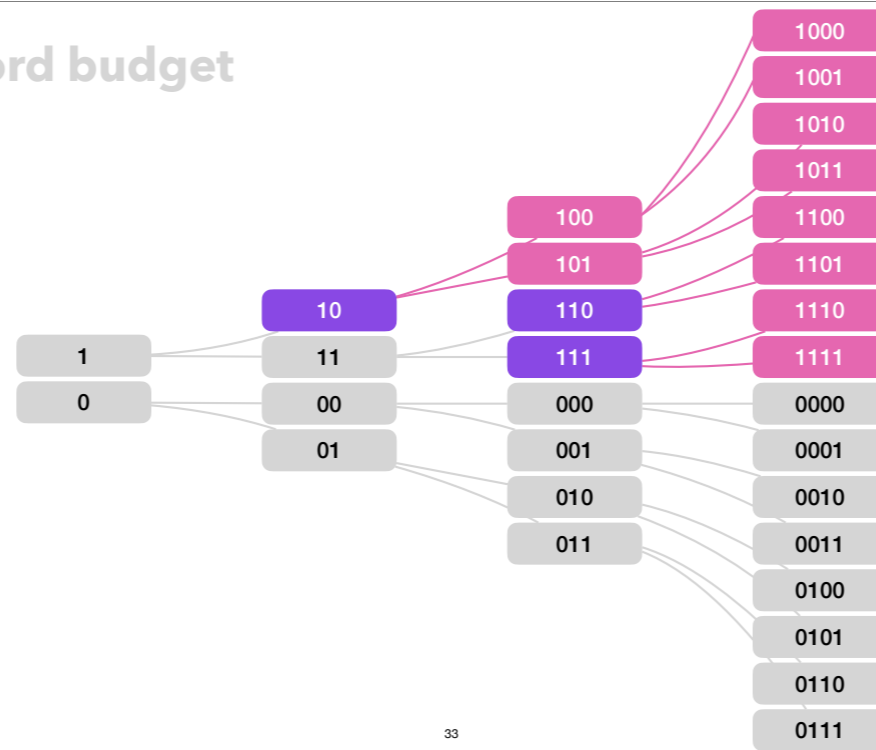$$2^{l_i}$$

this is assuming prefix-free codes, meaning that if we use e.g. '10' as a codeword, we can't use any other that begins the same way, otherwise after reading '10' we wouldn't know if the codeword was ending or a if we should continue reading

codeword budget

codeword budget

Kraft-McMillan inequality
proof: https://en.wikipedia.org/wiki/Kraft–McMillan_inequality#Proof_for_prefix_codes

# example coding function

$$p_1 = 1/16 \;\rightarrow\; l_1 = 4$$

| | | | |
|---|---|---|---|
| 1000 | 1100 | 0000 | 0100 |
| 1001 | 1101 | 0001 | 0101 |
| 1010 | 1110 | 0010 | 0110 |
| 1011 | 1111 | 0011 | 0111 |

- the grid is a probability distribution with 16 equally likely outcomes
- each element of the grid is given a unique codeword

$p_1 = 1/16 \;\rightarrow\; l_1 = 4$

| | | | |
|---|---|---|---|
| 1000 | 1100 | 0000 | 0100 |
| 1001 | 1101 | 0001 | 0101 |
| 1010 | 1110 | 0010 | 0110 |
| 1011 | 1111 | 0011 | 0111 |

# example coding function

$p_1 = 1/16 \;\rightarrow\; l_1 = 4$

| | | | |
|---|---|---|---|
| 1000 | 1100 | 0000 | 0100 |
| 1001 | 1101 | 0001 | 0101 |
| 1010 | 1110 | 0010 | 0110 |
|      | 1111 | 0011 | 0111 |

$p_3 = 1/8 \;\rightarrow\; l_3 = 3$

## example coding function

$p_1 = 1/16 \rightarrow l_1 = 4$

| | | | |
|---|---|---|---|
| 1000 | 1100 | 0000 | 0100 |
| 1001 | 1101 | 0001 | 0101 |
| 101 | 1110 | 0010 | 0110 |
| | 1111 | 0011 | 0111 |

$p_3 = 1/8 \rightarrow l_3 = 3$

38

if the two events in the bottom left become the same event with twice the probability, then substituting the new p_i=1/8 into the l_i=-log p_i equation will mean that we should find a length 3 codeword for it

$p_1 = 1/16 \;\rightarrow\; l_1 = 4$
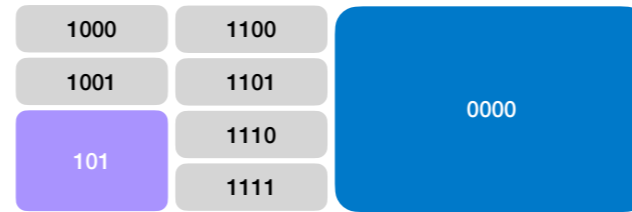
| | | | |
|---|---|---|---|
| 1000 | 1100 | 0000 | 0100 |
| 1001 | 1101 | 0001 | 0101 |
| 101 | 1110 | 0010 | 0110 |
| | 1111 | 0011 | 0111 |

$p_3 = 1/8 \;\rightarrow\; l_3 = 3$

$p_1 = 1/16 \;\rightarrow\; l_1 = 4$

| | |
|---|---|
| 1000 | 1100 |
| 1001 | 1101 |
| | 1110 |
| 101 | 1111 |

0000

$p_3 = 1/8 \;\rightarrow\; l_3 = 3$

$p_9 = 1/2 \;\rightarrow\; l_9 = 1$

$p_1 = 1/16 \;\rightarrow\; l_1 = 4$

| 1000 | 1100 | |
|------|------|---|
| 1001 | 1101 | 0 |
| 101 | 1110 | |
| | 1111 | |

$p_3 = 1/8 \;\rightarrow\; l_3 = 3$

$p_9 = 1/2 \;\rightarrow\; l_9 = 1$

$$p_i = 2^{-l_i}$$

probabilistic models

$$l_i = -\log p_i$$

possible codes (languages)

**find a model that describes the world well** ⟷ **find a 'language' in which experiences can be described concisely**
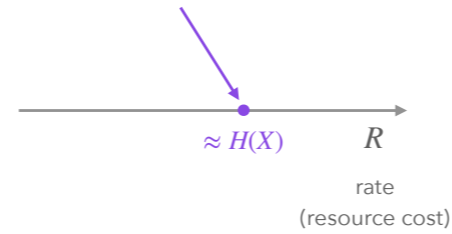
- it is possible to compress inputs (describe more concisely) without losing information
- the trick behind this is that if we know the frequencies/probability distribution of inputs, we can use short codewords for frequent inputs and on average we will then have a short description length
- there is a precise relationship between the frequencies and code lengths, see equations
- this establishes a correspondence between probabilistic models and encodings: in general if you have an encoding or description language, the description lengths imply a probability distribution, which is what the 'creator' of the language expects to happen

## lossless compression

(source coding theorem)

$$H(X) = \mathbb{E}[-\log p_i] \leq \mathbb{E}[l_i] < \mathbb{E}[-\log p_i] + 1$$
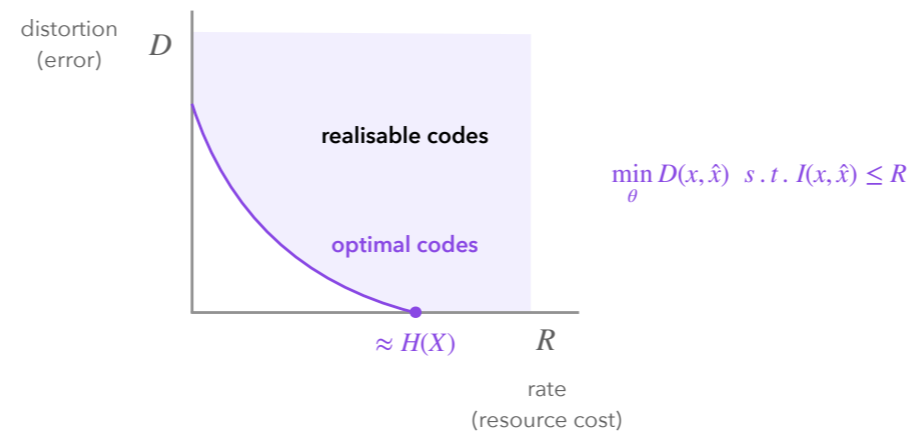
optimal lossless code

$\approx H(X)$     $R$

rate
(resource cost)

- the degree to which it is physically possible to compress without losing any details (lossless compression) is a property of the probability distribution of the input source (X). The numerical bound is given by the *entropy* H(X), which is the average of log(1/p_i).
- this quantity, H(X)=E[-log p_i], is the lower bound on the expected codeword length for an optimal lossless code
- for a detailed description of this and the tutorial's material, see MacKay Chapter 3 (p 67-81): http://www.inference.org.uk/itprnn/book.pdf

## lossy compression



distortion (error) — $D$

realisable codes

optimal codes

$\min_{\theta} D(x, \hat{x}) \ \ s.t. \ I(x, \hat{x}) \leq R$

$\approx H(X)$

$R$

rate (resource cost)

- if we want to compress the input even further, we will have to accept some loss of information, that is, some distortion in the reconstruction
- this means that we now have two axes along which we measure compression algorithms: the amount of memory resources similarly to as before, but now we also measure the expected distortion in the input

*we did not end up going into this so this is just for interest, not for exam:
rate here is measured in I(A,B), meaning mutual information between A and B
- this is the information that A contains about B (or that B contains about A)
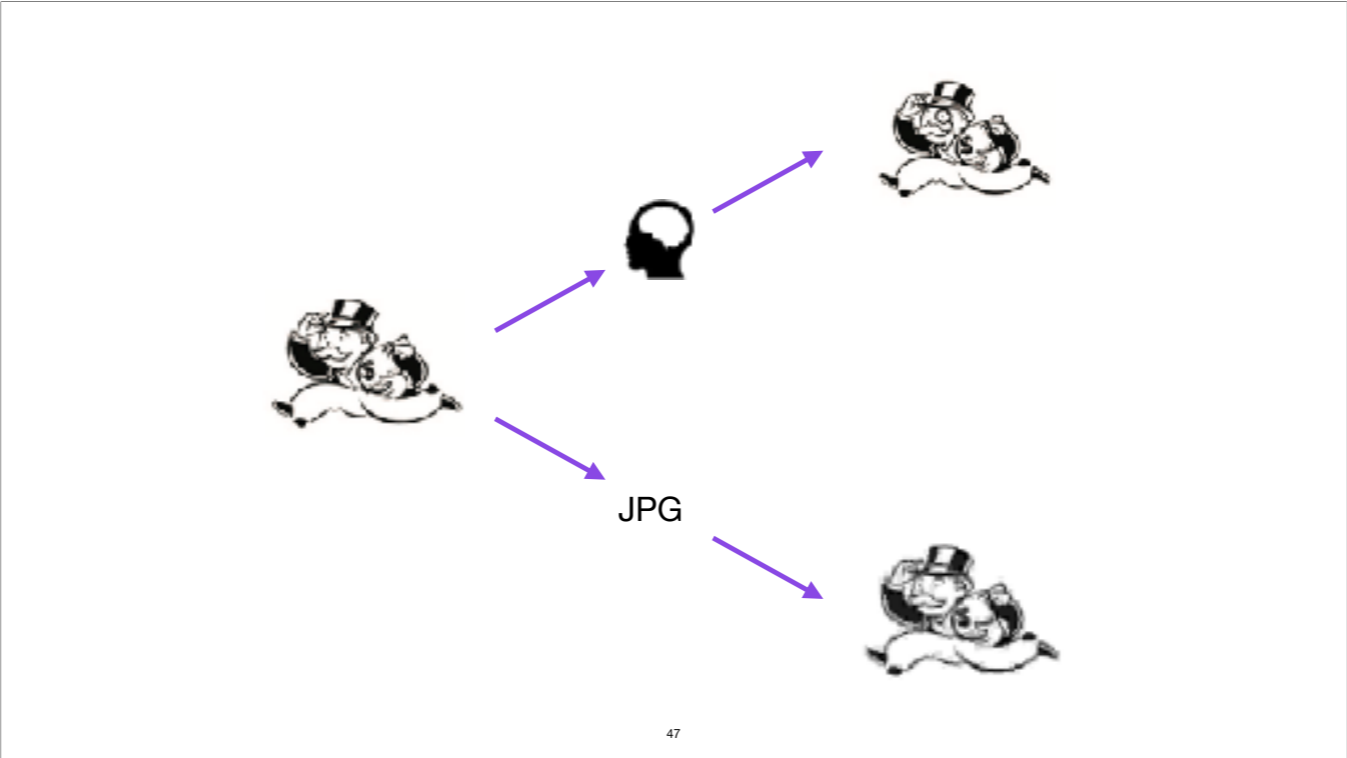- I(A,B)=H(A)-H(A|B)=H(B)-H(B|A)
- that is, the reduction in entropy about A if we learn the value of B
- if B is the memory trace for A, this tells us how noisy vs reliable the memory device is

# rate distortion trade-off



distortion
(error)

$D$

$\approx H(X)$    $R$

rate
(resource cost)

JPG

# generative compression



$$\hat{P}(x \mid \theta)$$
**generative model**

$$P(x)$$
**environment**

a generative model is a probabilistic model of how the environment generates observations

# generative compression

**semantic memory**
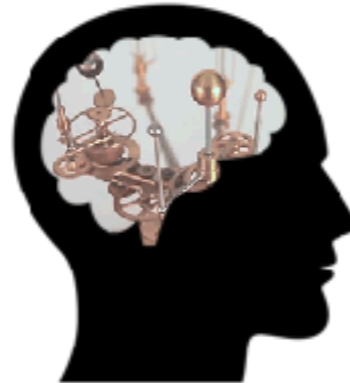
$$\hat{P}(x \mid \theta)$$

**sensory experience**

$$x$$

# perception as inference

**what are we interested in**

- what objects are around us
- how far
- who are around us
- what are they thinking
- what is going to happen

inference

**what can we observe**

- incoming photons
- air vibrations
- temperature fluctuations
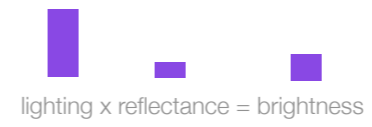- certain molecules





50

perception as inference

# perception as inference

snow in the evening seems white

perception as inference

54

coal in the sun still seems black, even though more photons arrive from it to our eyes than snow in the dark
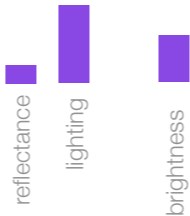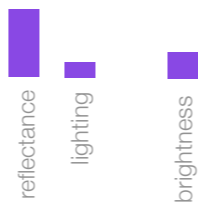
# perception as inference



lighting x reflectance = brightness

lighting

reflectance
(material)

brightness

*inference*

# perception as inference



lighting x reflectance = brightness

lighting

reflectance
(material)

brightness

inference

# perception as inference



lighting x reflectance = brightness

reflectance    lighting

brightness
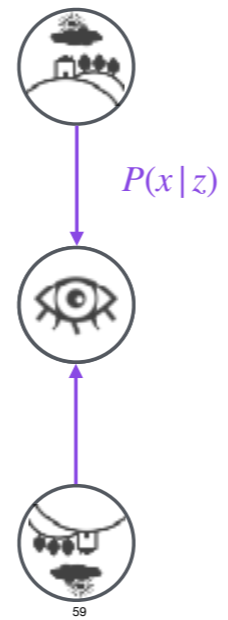
lighting      reflectance (material)

inference

57

# perception as inference

- generative direction

  - if the environment was in this state, what would I see?
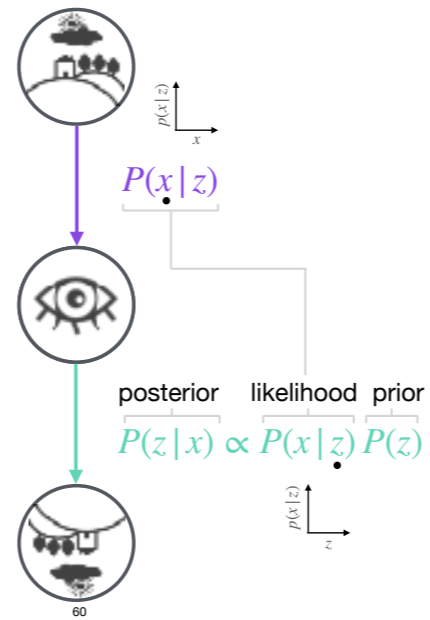
$$P(x \mid z)$$

**perception as inference**



$P(x\,|\,z)$

# perception as inference



- inverse direction

  - if I see this, what state is the environment in?

    - inverting the generative model

    - "vision is inverse graphics"

    - Bayesian inference

    - recognition model

$$P(x \mid z)$$

$$\underbrace{P(z \mid x)}_{\text{posterior}} \propto \underbrace{P(x \mid z)}_{\text{likelihood}} \underbrace{P(z)}_{\text{prior}}$$

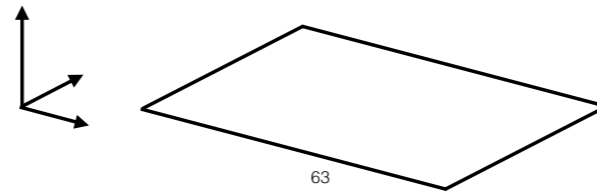60

# perception as inference
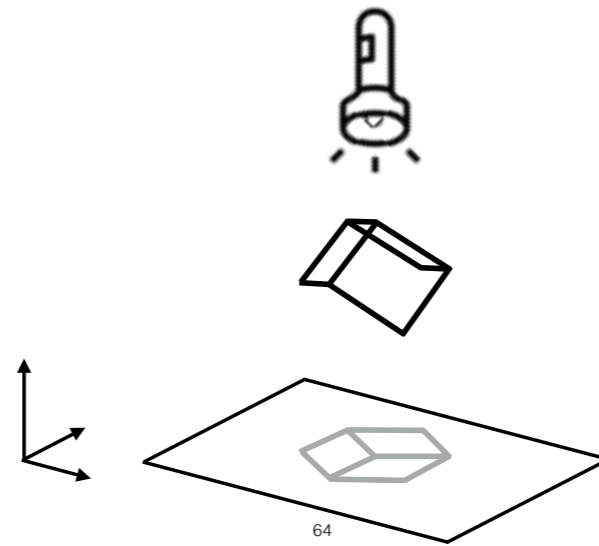
# perception as inference



- what we see is not the data but an interpretation of the data
- 'unconscious inference' over latent variables
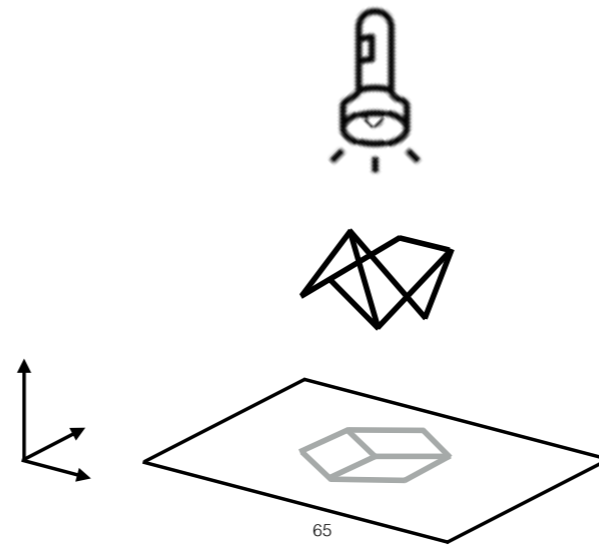
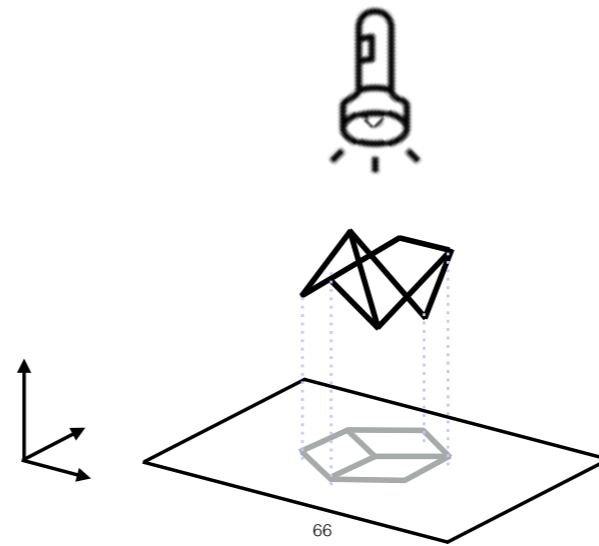# perception as inference

Kersten & Yuille, 2003

64

Kersten & Yuille, 2003

# perception as inference

Kersten & Yuille, 2003

# perception as inference

Kersten & Yuille, 2003

**perception as inference**

Kersten & Yuille, 2003
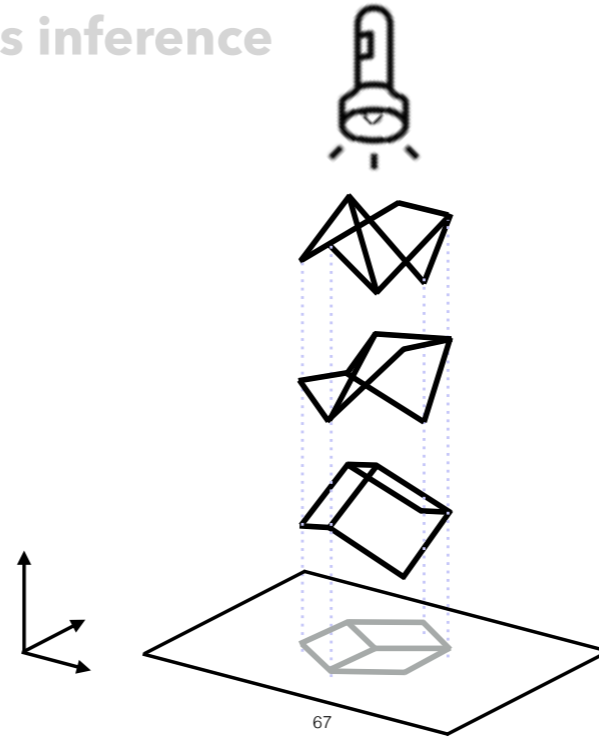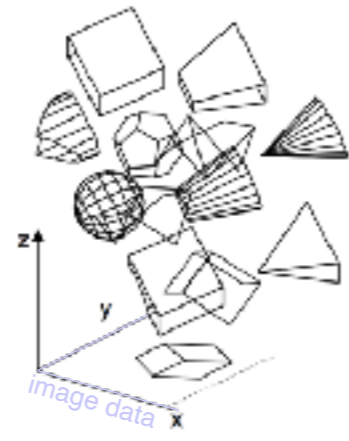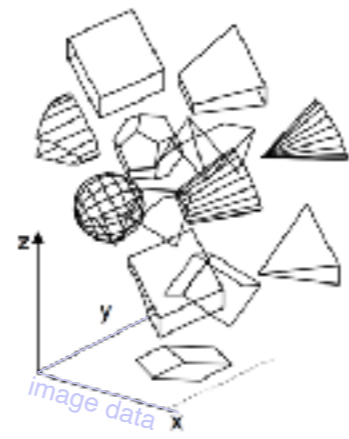
**perception as inference**
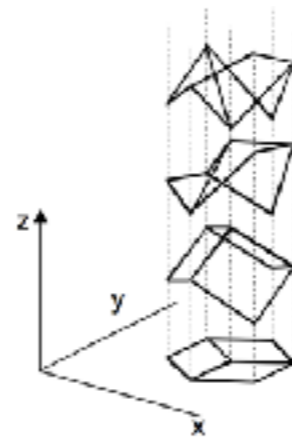


hypotheses sampled
from prior

68

# perception as inference



hypotheses sampled
from prior

hypotheses with
nonzero likelihood

Kersten & Yuille, 2003

# perception as inference



hypotheses sampled
from prior

hypotheses with
nonzero likelihood

Kersten & Yuille, 2003

# perception as inference

"the girl saw the boy with the telescope"



$j_1$                    $j_2$

71

Goodman & Frank, 2016

**perception as inference**



$P(m \mid u)$

$u$

72

**perception as inference**

McGurk & MacDonald, 1976

**variational autoencoder**

$x$

input

$z$

memory trace
(code)

$\hat{x}$

reconstruction

Kingma & Welling, 2014; Rezende et al, 2014
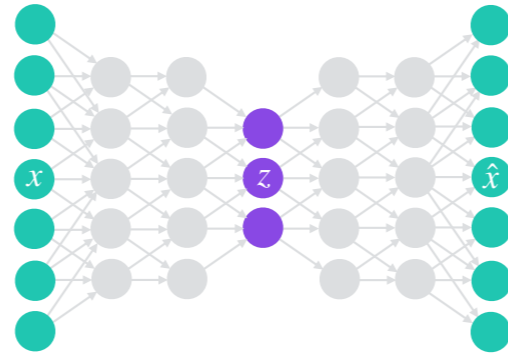
beta-VAE can be seen as both as a generative model + an inference method for approximately inverting it, and also as a lossy compression algorithm

variational autoencoder

variational autoencoder

inference

prior

generate $\hat{x}$

$\hat{x} \sim p(x \mid z)$

sample z

$z \sim q(z \mid x)$

$z$

$x$

compute approximate posterior

$Z$

generation

$\hat{x}$

$q(z \mid x) \approx p(z \mid x)$

Kingma & Welling, 2014; Rezende et al, 2014

1. compute an approximate posterior based on sample x
2. sample a single z from this posterior
3. conditioning the generative network on this sample z, generate a reconstruction

variational autoencoder
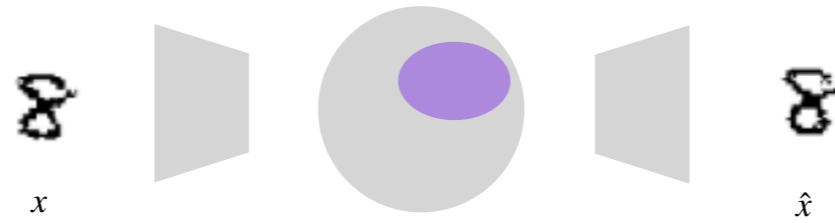
$x$

$\hat{x}$

Kingma & Welling, 2014; Rezende et al, 2014

different x-es will correspond to different posteriors and therefore different reconstructions, but in case the posteriors overlap, the memory traces for different stimuli might be confused

# variational autoencoder



$x$

$\hat{x}$

Kingma & Welling, 2014; Rezende et al, 2014

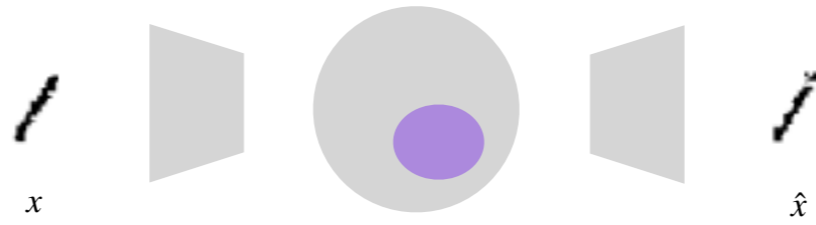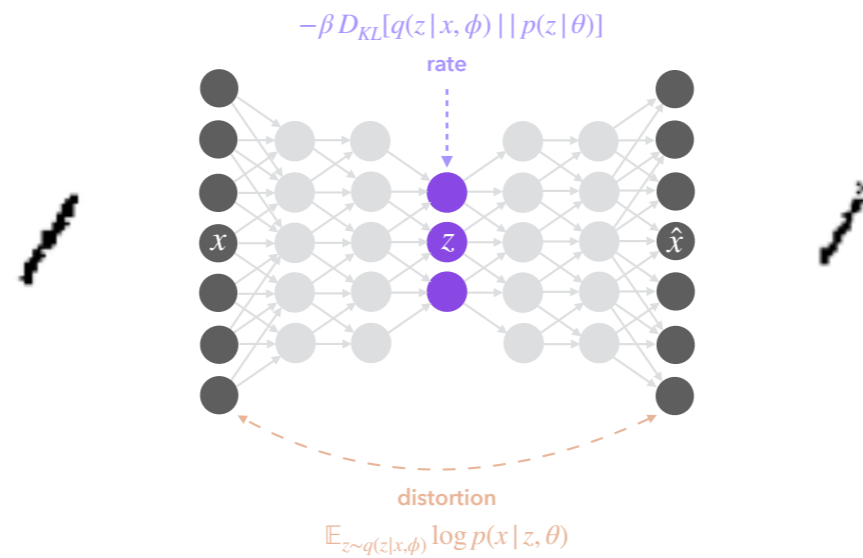# variational autoencoder



$x$

$\hat{x}$

Kingma & Welling, 2014; Rezende et al, 2014

variational autoencoder

$$-\beta D_{KL}[q(z\,|\,x,\phi)\,||\,p(z\,|\,\theta)]$$

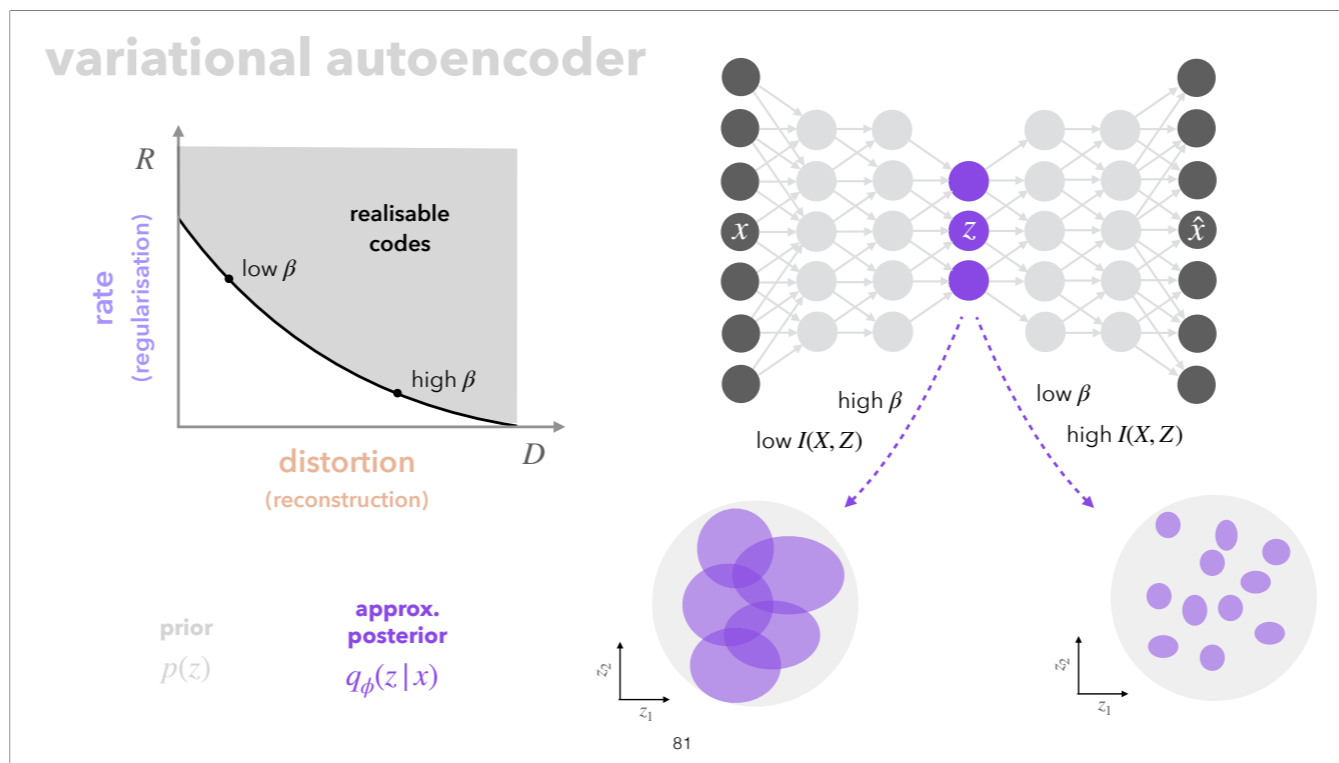rate

$z$

distortion

$$\mathbb{E}_{z\sim q(z|x,\phi)}\log p(x\,|\,z,\theta)$$

80

Kingma & Welling, 2014; Rezende et al, 2014

- in a VAE, both the inference method and the generative model are neural networks
- the objective function for training the beta-VAE consists of two terms, which in the compression view correspond to rate and distortion. In the generative modelling literature these are called 'regularisation term' and 'reconstruction' term respectively
- D_KL is KL-divergence, a form of (almost) distance between probability distributions

variational autoencoder

realisable codes

low $\beta$

high $\beta$

$R$ — rate (regularisation)

$D$ — distortion (reconstruction)

high $\beta$
low $I(X, Z)$

low $\beta$
high $I(X, Z)$

prior
$p(z)$

approx. posterior
$q_\phi(z \mid x)$

$x$ $z$ $\hat{x}$

$z_1$ $z_2$

81

- intuitively, when KL term is weighted more strongly (using beta), posteriors have to resemble the wide and spherical prior more
- this makes it more likely that they overlap and makes the variance of the sampled z larger
- this corresponds to a lower rate (less information in the memory trace about the original input)
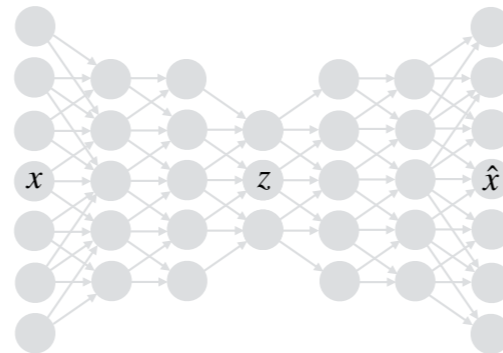
## variational autoencoder

**chess positions**

**sketches**

$x$  $z$  $\hat{x}$

**word lists**

butter
food
eat
sandwich
…

$x$

3000 chess games
(FICS database)

quickdraw75k
object pairs

food
butter
bread
sandwich
…

(small subset of)
wikipedia

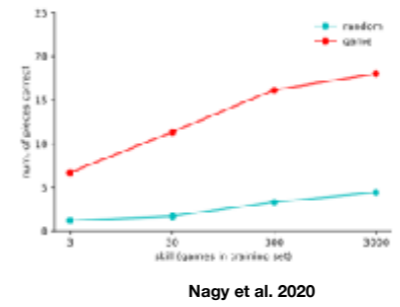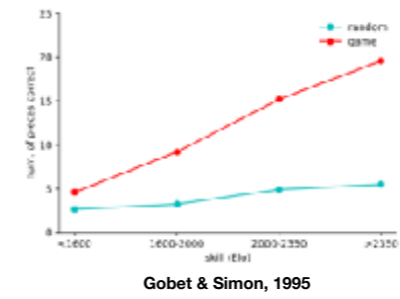$\sim P(\hat{x}\,|\,x)$  $\mathscr{D}$

Nagy et al. 2020
Bates & Jacobs 2020, Hedayati et al.

82

# domain expertise

reconstruct state of chess board after 5 seconds viewing time

domain expertise

Gobet & Simon, 1995

Nagy et al. 2020

84

- accuracy increases with domain expertise, but only for configurations from actual chess games
- for randomly shuffled configurations, skill does not help much
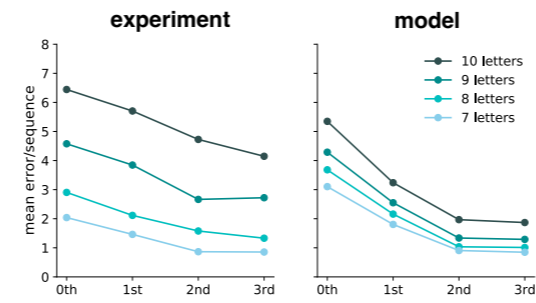- skill in model is amount of training of beta-VAE

# domain congruence

| uniform | 1st order | 2nd order | 3rd order | full word |
|---------|-----------|-----------|-----------|-----------|
| RCIFODWVIL | TNEOOESHHE | HIRTOCLENO | BETEREASYS | PLANTATION |
| GKTODKPENF | INOLGGOLVN | DOVEECOFOF | CRAGETTERS | FLASHLIGHT |
| TZXKHAWCCF | PDOASLOTPP | SESERAICCG | TOWERSIBLE | UNCOMMONLY |
| NGORHQIYWB | AEOCAOIAON | AREDAGORTZ | DEEMEREANY | ALIENATION |
| BVNJSYZXUA | IRCRENFCTN | CUNSIGOSUR | THERSERCHE | PICKPOCKET |

85

- reconstruct words with statistics that are increasingly congruent with the statistics of english language
- very similar to chess example, but there we had two degrees of congruence - random and game
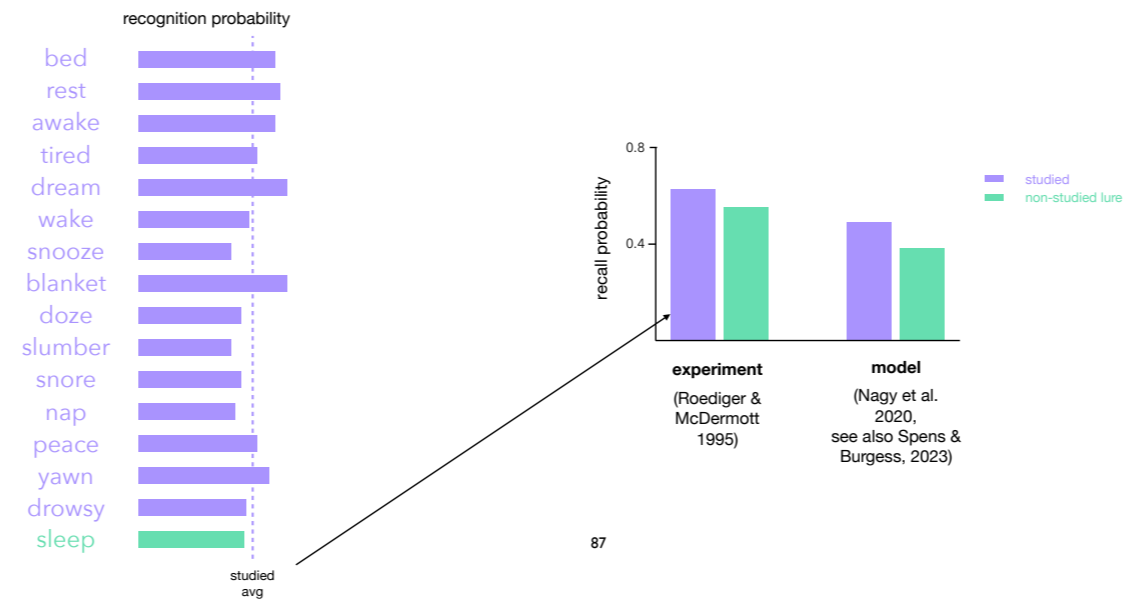
## domain congruence

**experiment**

**model**

mean error/sequence

- 10 letters
- 9 letters
- 8 letters
- 7 letters

0th   1st   2nd   3rd      0th   1st   2nd   3rd

| uniform | 1st order | 2nd order | 3rd order | full word |
|---------|-----------|-----------|-----------|-----------|
| RCIFODWVIL | TNEOOESHHE | HIRTOCLENO | BETEREASYS | PLANTATION |
| GKTODKPENF | INOLGGOLVN | DOVEECOFOF | CRAGETTERS | FLASHLIGHT |
| TZXKHAWCCF | PDOASLOTPP | SESERAICCG | TOWERSIBLE | UNCOMMONLY |
| NGORHQIYWB | AEOCAOIAON | AREDAGORTZ | DEEMEREANY | ALIENATION |
| BVNJSYZXUA | IRCRENFCTN | CUNSIGOSUR | THERSERCHE | PICKPOCKET |

(experiment: Baddeley et al., 1971, model: Frater et al., 2022)

86

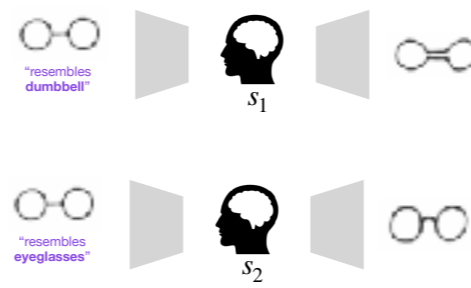- accuracy increases with degree of congruence with english, and decreases with length of word

**gist-based distortions**

recognition probability

bed
rest
awake
tired
dream
wake
snooze
blanket
doze
slumber
snore
nap
peace
yawn
drowsy
sleep

studied
avg

recall probability

0.8

0.4

studied
non-studied lure

**experiment**
(Roediger &
McDermott
1995)

**model**
(Nagy et al.
2020,
see also Spens &
Burgess, 2023)

87

– lure is a semantically related word that was not in the list
– lures are falsely recognised with comparable probability to studied items
– effect of regenerating the list from stored latent representation

gist-based distortions

"resembles dumbbell"

$s_1$

"resembles eyeglasses"

$s_2$

88

- subjects view ambiguous stimuli, with different labels for different subjects
- labels introduce label-specific distortions in recall
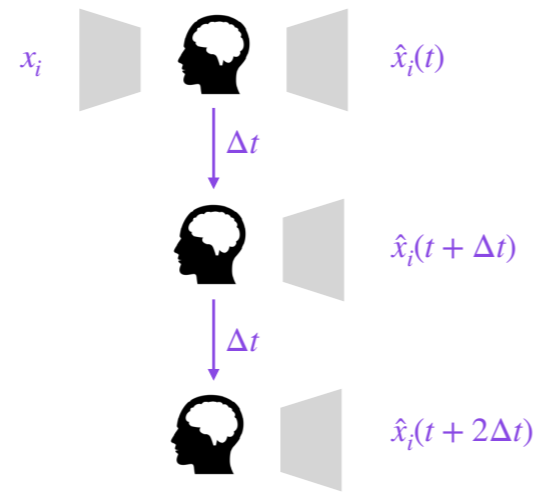
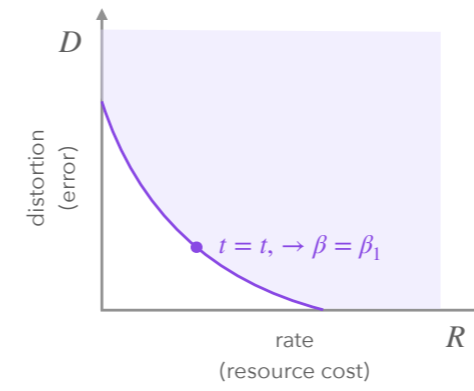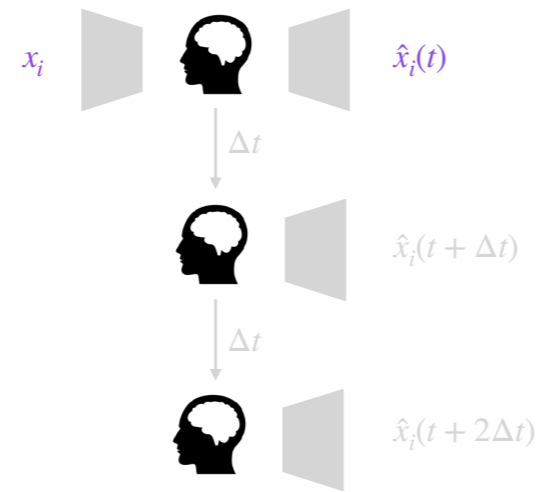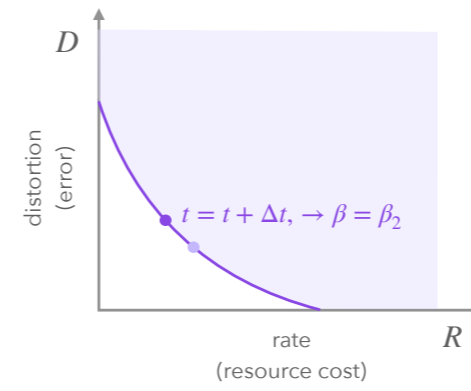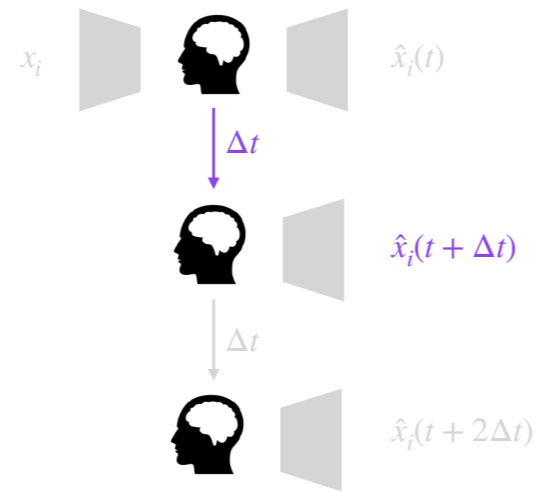# gist-based distortions



Carmichael et al., 1932

# gist-based distortions



Carmichael et al., 1932

Nagy et al. 2020

# delayed recall

$x_i$  $\hat{x}_i(t)$

$\Delta t$

 $\hat{x}_i(t + \Delta t)$

$\Delta t$

 $\hat{x}_i(t + 2\Delta t)$

$$\Delta t$$

$x_1 \quad x_2 \quad x_3 \quad x_2 \quad x_1 \quad x_4$



accuracy

$\Delta t$

# delayed recall

$x_i$ $\hat{x}_i(t)$

$\Delta t$

$\hat{x}_i(t + \Delta t)$

$\Delta t$

$\hat{x}_i(t + 2\Delta t)$

$D$

distortion
(error)

$t = t, \rightarrow \beta = \beta_1$

rate
(resource cost)

$R$

**delayed recall**

$x_i$

$\hat{x}_i(t)$

$\Delta t$

$\hat{x}_i(t + \Delta t)$

$\Delta t$

$\hat{x}_i(t + 2\Delta t)$

$D$

distortion (error)

$t = t + \Delta t, \rightarrow \beta = \beta_2$

rate (resource cost)

$R$

**delayed recall**

$x_i$  $\hat{x}_i(t)$

$\Delta t$

$\hat{x}_i(t + \Delta t)$

$\Delta t$

$\hat{x}_i(t + 2\Delta t)$



$D$

distortion (error)

$t = t + 2\Delta t, \rightarrow \beta = \beta_3$

rate (resource cost)

$R$

Nagy et al. 2020

# goals, tasks, value

agent

perception

action

environment

goals, tasks, value

$x_t$      $\hat{x}_t$
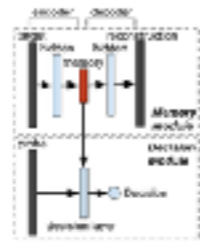
$$d(x, \hat{x}) = (x - x)^2$$

see Bates & Jacobs, 2020

- in engineering contexts, typical distortion is MSE in pixel space
- for a robot that needs to manipulate a small white ball, removing this ball from the input leads to almost negligible reconstruction error
- need to overweight errors that are relevant to rewards
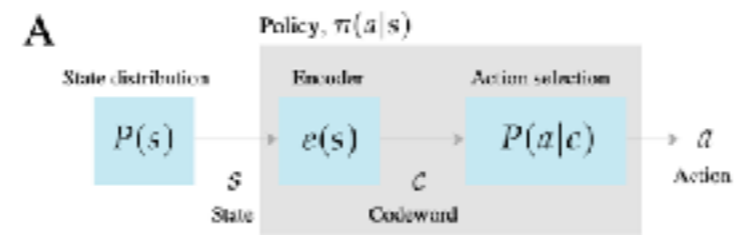
goals, tasks, value

$r_t$ $a_t$

$x_t$ $\hat{x}_t$

based on Hafner et al., 2023

98

- this can be incorporated in VAEs, for example this is the basis of the DREAMER v3 model that you've seen in lecture 5

* goals, tasks, value
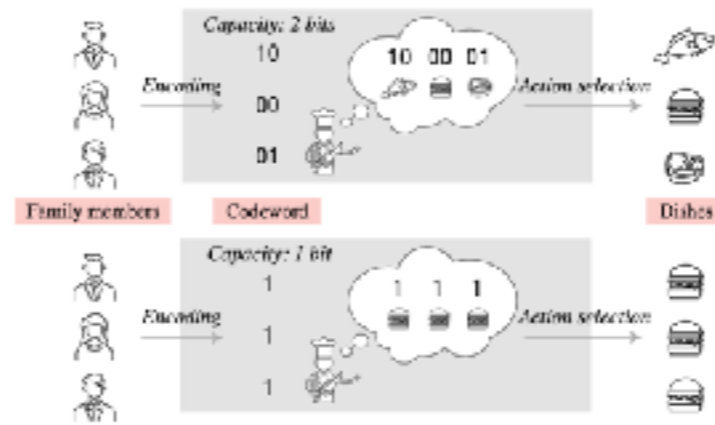
Bates & Jacobs, 2020

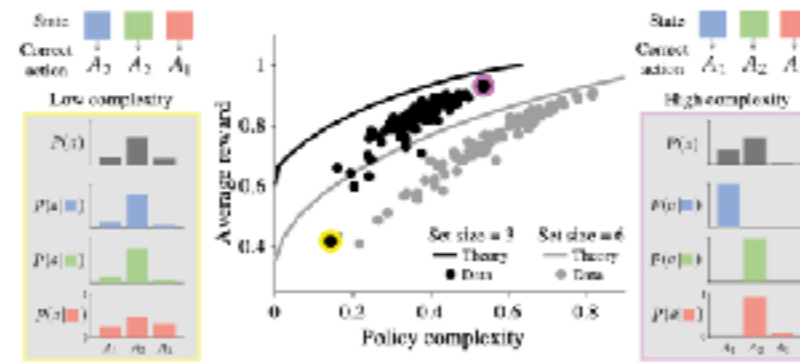- humans are also sensitive to reward-relevance in memory accuracy
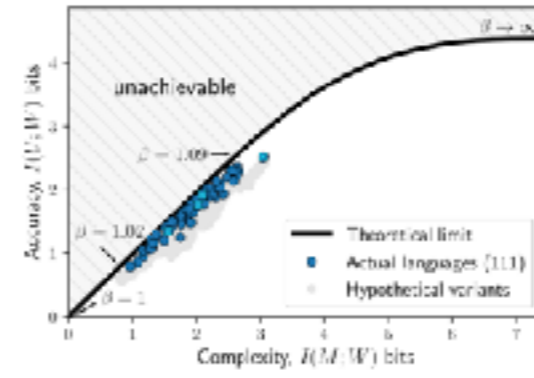
# * goals, tasks, value



**A**

State distribution

$P(s)$

Policy, $\pi(a|s)$

Encoder

$e(s)$

Action selection

$P(a|c)$

$a$

Action

$s$
State

$c$
Codeword

Lai & Gershman, 2021

# * goals, tasks, value

Lai & Gershman, 2021

# * goals, tasks, value

Lai & Gershman, 2021

**\***

Zaslavsky et al 2018

*

104

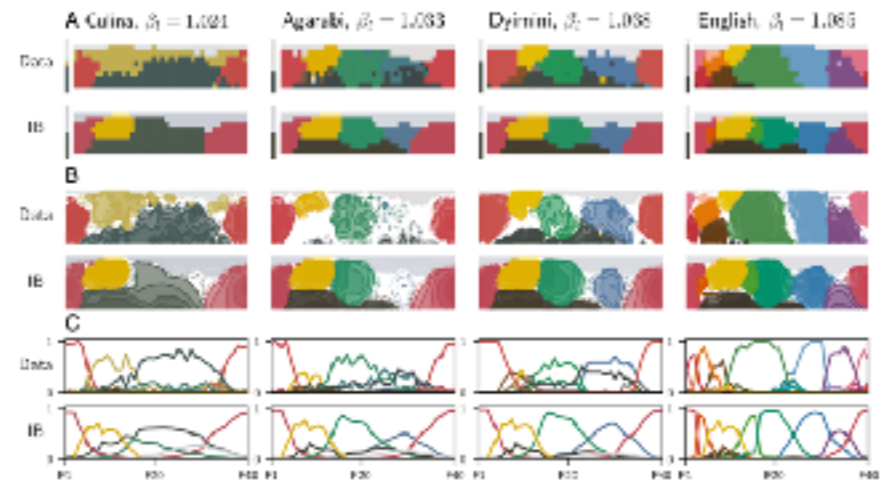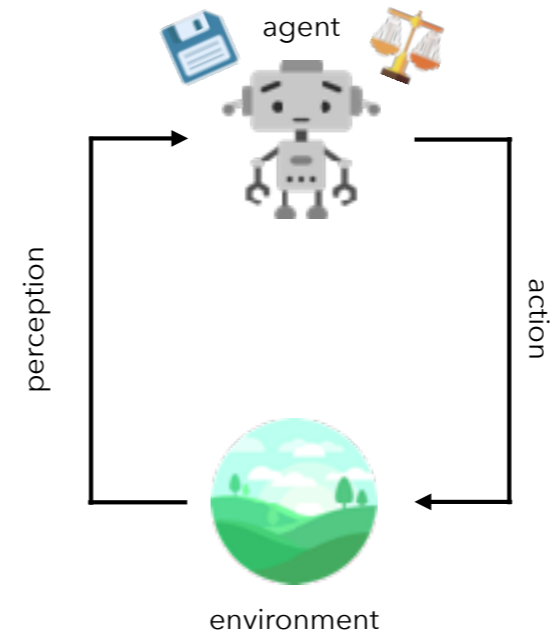## summary

- focus on constraints on **memory**

  - lossless compression

  - lossy compression

  - generative compression

    - perception as bayesian inference

    - human memory distortions

agent

perception

action

environment

## can information capture all constraints?

xkcd

Next week: Concepts & Categories

Is a hotdog a sandwich?