# Attention models & interpretability

Adrian Valente
30-11-2023
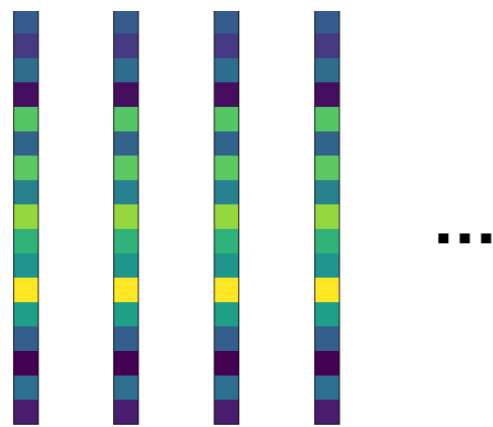
# Deep learning primer



0.12   .4
1.3   1.2
.9   2.5

"dog"

...   "cat"

"croissant"

**Goal: learn mapping vector -> vector**

# Problem: variable-length input

- Text

- Sound

- Video

- Time series

- ....

$$(\mathbf{x_1}, \mathbf{x_2}, \cdots)$$

# RNNs

**Example: sentiment analysis**

**Mapping sequence -> vector**



Life            is            really            tough
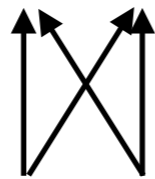
# RNNs

**Example: sentiment analysis**

**Mapping sequence -> vector**



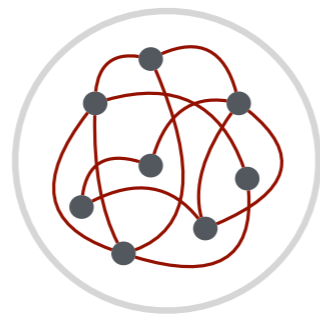Life         is        really        tough

# RNNs

**Example: sentiment analysis**

**Mapping sequence -> vector**



Life        is        really        tough

# RNNs

**Example: sentiment analysis**

**Mapping sequence -> vector**



positive

negative

is          really          tough

# Backprop through time

**Unrolled computation graph**



$\mathbf{h}_0$  $\mathbf{W}_{rec}$  $\mathbf{h}_1$  $\mathbf{W}_{rec}$  $\mathbf{h}_2$  $\mathbf{W}_{rec}$  $\mathbf{h}_3$  $\mathbf{W}_{rec}$  $\mathbf{h}_4$  $\mathbf{W}_{out}$  $\mathbf{y}$

$\mathbf{W}_{in}$  $\mathbf{W}_{in}$  $\mathbf{W}_{in}$  $\mathbf{W}_{in}$

Life            is            really            tough

# Backprop through time

**Unrolled computation graph**



$\mathbf{h}_0$ $\quad$ $\mathbf{W}_{rec}$ $\quad$ $\mathbf{h}_1$ $\quad$ $\mathbf{W}_{rec}$ $\quad$ $\mathbf{h}_2$ $\quad$ $\mathbf{W}_{rec}$ $\quad$ $\mathbf{h}_3$ $\quad$ $\mathbf{W}_{rec}$ $\quad$ $\mathbf{h}_4$ $\quad$ $\mathbf{W}_{out}$ $\quad$ $\mathbf{y}$

$\mathbf{W}_{in}$ $\qquad$ $\mathbf{W}_{in}$ $\qquad$ $\mathbf{W}_{in}$ $\qquad$ $\mathbf{W}_{in}$

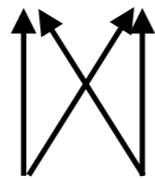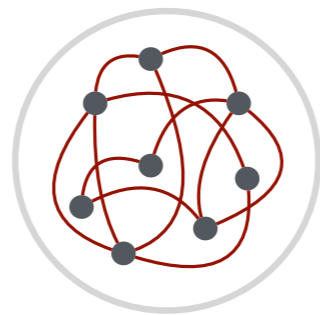Life $\qquad\qquad$ is $\qquad\qquad$ really $\qquad\qquad$ tough

# seq2seq (Sutskever et al. 2014)

**Example: language translation**

**Mapping sequence -> sequence**



My          father's          green          hat

# seq2seq (Sutskever et al. 2014)



My       father's       green       hat

# seq2seq (Sutskever et al. 2014)



My        father's        green        hat

# seq2seq (Sutskever et al. 2014)



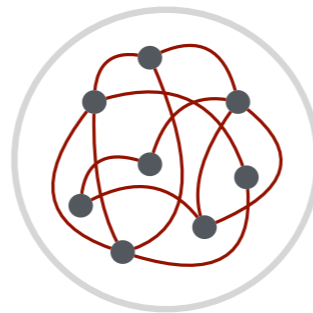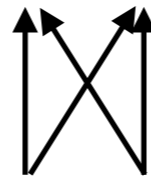father's          green          hat

# seq2seq (Sutskever et al. 2014)

Le



$\mathbf{h}_T$

# seq2seq (Sutskever et al. 2014)
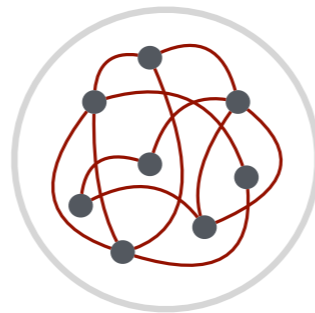
Le    chapeau

$$\mathbf{h}_{T+1}$$

# seq2seq (Sutskever et al. 2014)

Le    chapeau    vert

$\mathbf{h}_{T+2}$

# seq2seq (Sutskever et al. 2014)

Le    chapeau    vert    de
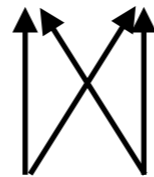


$\mathbf{h}_{T+3}$

# seq2seq (Sutskever et al. 2014)

Le    chapeau    vert    de    mon

$\mathbf{h}_{T+4}$

# seq2seq (Sutskever et al. 2014)

Le   chapeau   vert   de   mon   père



$\mathbf{h}_{T+5}$

# seq2seq (Sutskever et al. 2014)

**Unrolled graph**

# Non-linear processing of the sequence

My father's green hat

Le chapeau vert de mon père

# Non-linear processing of the sequence

My father's green hat

Le chapeau vert de mon père

# Non-linear processing of the sequence

My father's green hat

Le chapeau vert de mon père

# Attention

*"It is the taking possession by the mind […] of one out of what seem several simultaneously possible objects or trains of thought."*

William James, 1890

Modern analogy:
It is a "spotlight"

# An artificial attentional mechanism for translation

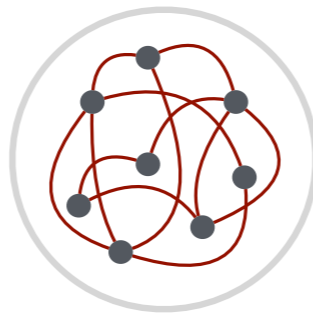Bahdanau-Cho-Bengio, *Neural machine translation by jointly learning to align and translate, ICLR 2015*

# An artificial attentional mechanism for translation

$$\mathbf{c}_t = \sum_j \alpha_{tj}\mathbf{x}_j$$

with

$$\sum_j \alpha_{tj} = 1$$

$$\alpha_{tj} \in [0,1]$$

$\alpha_{t1}$

$\alpha_{t2}$    $\alpha_{t3}$

$\alpha_{t4}$

My          father's          green          hat

# An artificial attentional mechanism for translation

Le        chapeau

*output network*

$$\mathbf{c}_t = \sum_j \alpha_{tj} \mathbf{x}_j$$

with

$$\sum_j \alpha_{tj} = 1$$

$$\alpha_{tj} \in [0,1]$$

My        father's        green        hat

# An artificial attentional mechanism for translation

*controller RNN*

Le      chapeau

*output network*

$$\mathbf{c}_t = \sum_j \alpha_{tj} \mathbf{x}_j$$

with

$$\sum_j \alpha_{tj} = 1$$

$$\alpha_{tj} \in [0, 1]$$

$\{\alpha_{t+1,j}\}$

+

My      father's      green      hat

# An artificial attentional mechanism for translation



*Inherently opens a window on the model's internal process!!*

# Visual attention for image captioning

Xu et al., *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*, 2015



A woman is throwing a <u>frisbee</u> in a park.

A little <u>girl</u> sitting on a bed with a teddy bear.

# Transformers



[Vaswani et al 2017]

# Transformers: self-attention



My          father's          green          hat

# Transformers: self-attention

# Transformers: working definition of attention

*"An attention function can be described as mapping a query and a set of key-value pairs to an output [...] The output is computed as a weighted sum of the values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key."*



$$\alpha_{i,j} = \mathrm{softmax}(\mathbf{q}_i^T \mathbf{k}_i)\mathbf{v}_i$$

# Transformers: working definition of attention

*"An attention function can be described as mapping a query and a set of key-value pairs to an output […] The output is computed as a weighted sum of the values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key."*

$$\alpha_{i,j} = \operatorname{softmax}(\mathbf{q}_i^T \mathbf{k}_i)\mathbf{v}_i$$

output 1

Query 1 ——→ **Attention**

Query 2

$\alpha_1$

$\alpha_2$

$\alpha_3$

value 1,
key1

value 2,
key2

value 3,
key 3
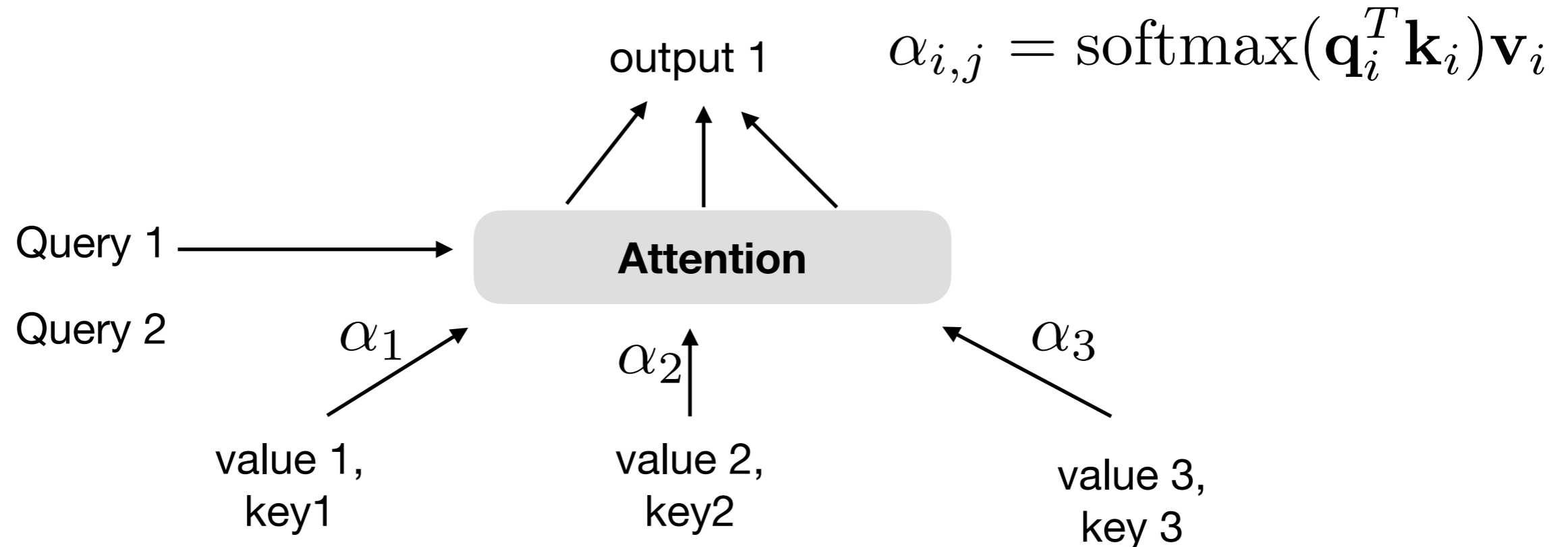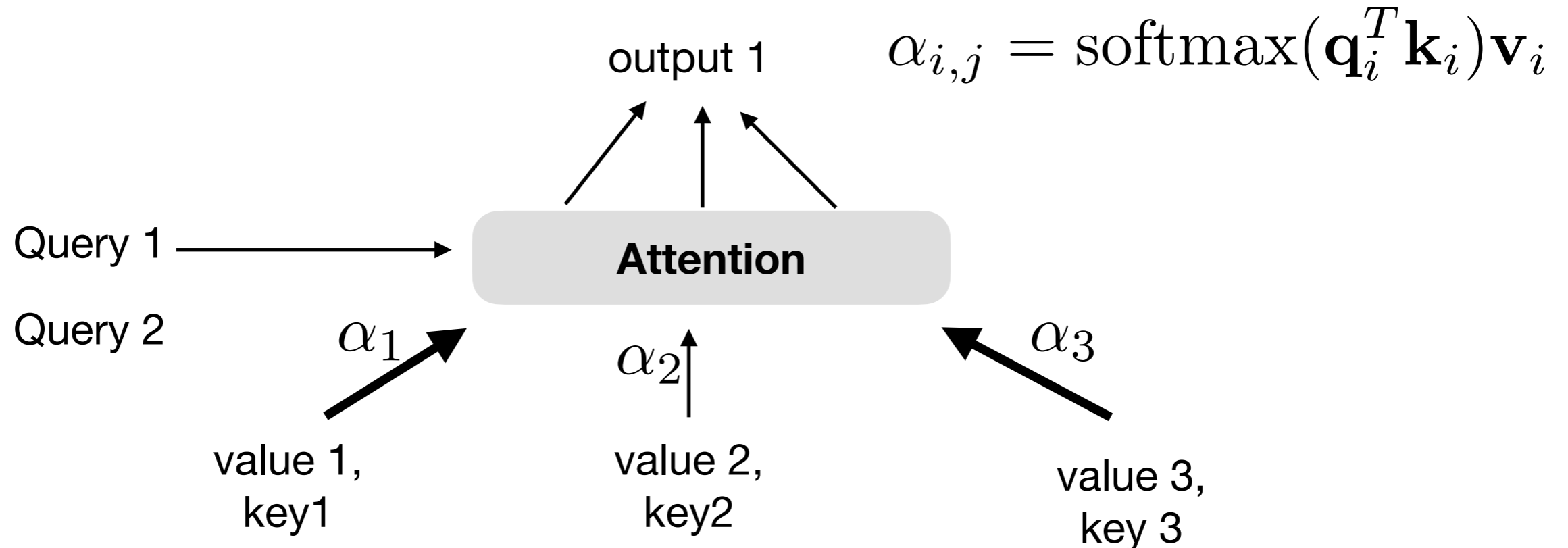
# Transformers: working definition of attention

*"An attention function can be described as mapping a query and a set of key-value pairs to an output [...] The output is computed as a weighted sum of the values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key."*
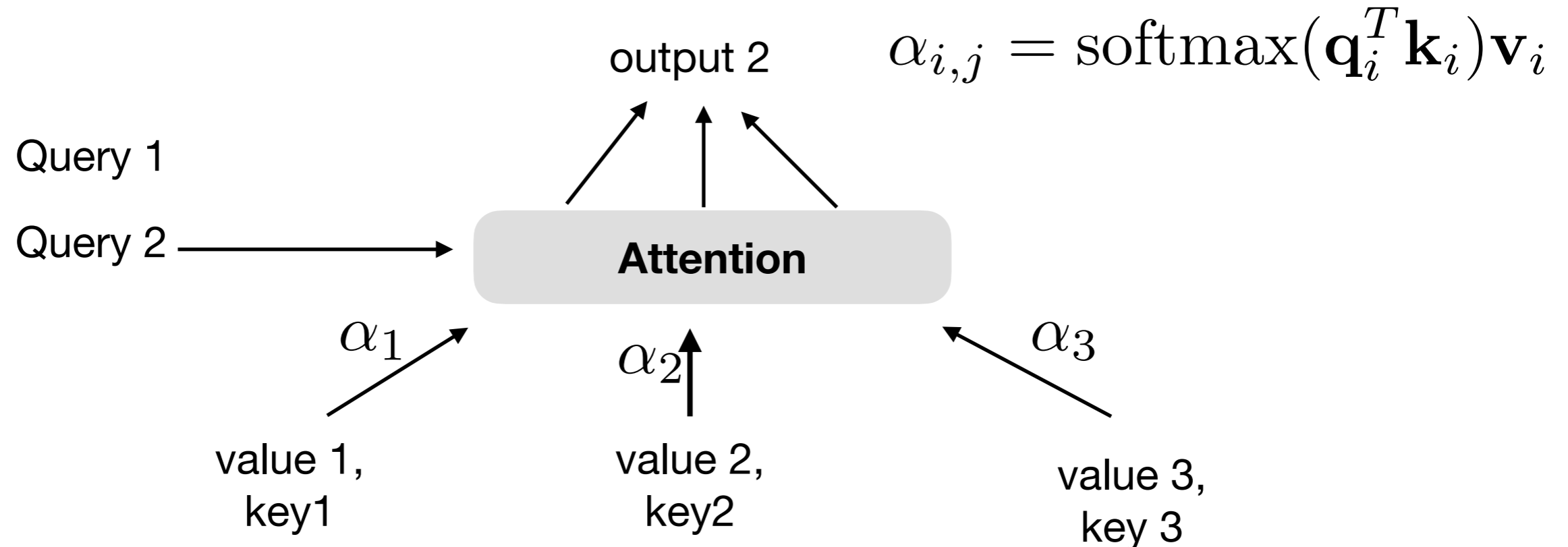
output 2

$$\alpha_{i,j} = \text{softmax}(\mathbf{q}_i^T \mathbf{k}_i)\mathbf{v}_i$$

Query 1

Query 2 $\longrightarrow$ **Attention**

$\alpha_1$  $\alpha_2$  $\alpha_3$

value 1,
key1

value 2,
key2

value 3,
key 3

# Transformers: working definition of attention

Input    Thinking    Machines

Embedding   $X_1$    $X_2$

Queries   $q_1$    $q_2$    $W^Q$    $\mathbf{q}_i = W_Q \mathbf{x}_i$

Keys   $k_1$    $k_2$    $W^K$    $\mathbf{k}_i = W_K \mathbf{x}_i$

Values   $v_1$    $v_2$    $W^V$    $\mathbf{v}_i = W_V \mathbf{x}_i$

# In matrix form

$$\text{softmax}\left(\frac{\boxed{Q} \times \boxed{K^T}}{\sqrt{d_k}}\right) \boxed{V}$$

$$= \boxed{Z}$$

# Transformers: self-attention



[Jesse Vig, 2019]

# Interpretability through attention



Figure 4: Attention pattern in GPT-2 related to coreference resolution suggests the model may encode gender bias.

[Jesse Vig, 2019]

# Self-attention can process variable-length inputs



$$a(\cdot, \{\dots\}) : \mathbf{x}_i, \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \to \mathbf{y}_i$$

# Gritty details: multi-head attention

# Gritty details: multi-head attention



other issue: information can get mixed across tokens layer by layer, and attention doesn't necessarily represent attention to the corresponding word anymore (see Darcet et al 2023 for example).

# Gritty details: encoder-decoders

Encoder: e.g. BERT
Typical task: Masked language modelling

# Gritty details: encoder-decoders

Decoder: e.g. GPT
Typical task: Autoregressive language modelling

# Gritty details: encoder-decoders

Encoder-decoder: e.g. T5
Typical task: Sequence-to-sequence modelling

# Gritty details: encoder-decoders

Decoding time step: (1) 2 3 4 5 6                    OUTPUT

# Gritty details: encoder-decoders

Decoding time step: 1 ②3 4 5 6

OUTPUT I

K encdec   V encdec

Linear + Softmax

ENCODERS

DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT   Je   suis   étudiant

PREVIOUS OUTPUTS   I

# Gritty details: positional embeddings



$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

# FF layers & attention

# How to make sense?
# Embeddings



Male-Female          Verb Tense          Country-Capital

**Semantic geometrical encoding**

[eg. Mikolov et al. 2013]

# How to make sense? Embeddings



[Lample et al. 2018, Word translation without parallel data]

# Context-dependent embeddings (BERT)



German article "die"

Was der Fall ist, **die** Tatsache, ist das Bestehen von Sachverhalten.

über **die** Verhandlungen der Königl.

single person dies ←——→ multiple people die

Chernenko became the first Soviet leader to **die** in less than three years

Vaughan's ultimate fantasy was to **die** in a head-on collision with movie star Elizabeth Taylor

Over 60 people **die** and over 100 are unaccounted for.

Many more **die** from radiation sickness, starvation and cold.

a playing die

Players must always move a token according to the **die** value

The faces of a **die** may be placed clockwise or counterclockwise

**Semantic geometrical encoding**

[Coenen et al., 2019]

# FF layers as key-value stores



$$p(k_i \mid x) \propto \exp(\mathbf{x} \cdot \mathbf{k}_i)$$

$$\mathbf{MN}(\mathbf{x}) = \sum_{i=1}^{d_m} p(k_i \mid x)\mathbf{v}_i$$

[Geva et al., 2021]

# Causal scrubbing



(a) Clean run — The, Space, Need, le, is, in, downtown → Seattle (correct output)

(b) Corrupted subject run — The*, Space*, Need*, le*, is, in, downtown → ? (corrupted output)

(c) Patch clean states

(d) Note when output is fixed

Legend: $h_i^{(l)}$ state; attention; MLP; corrupted embedding; example flow

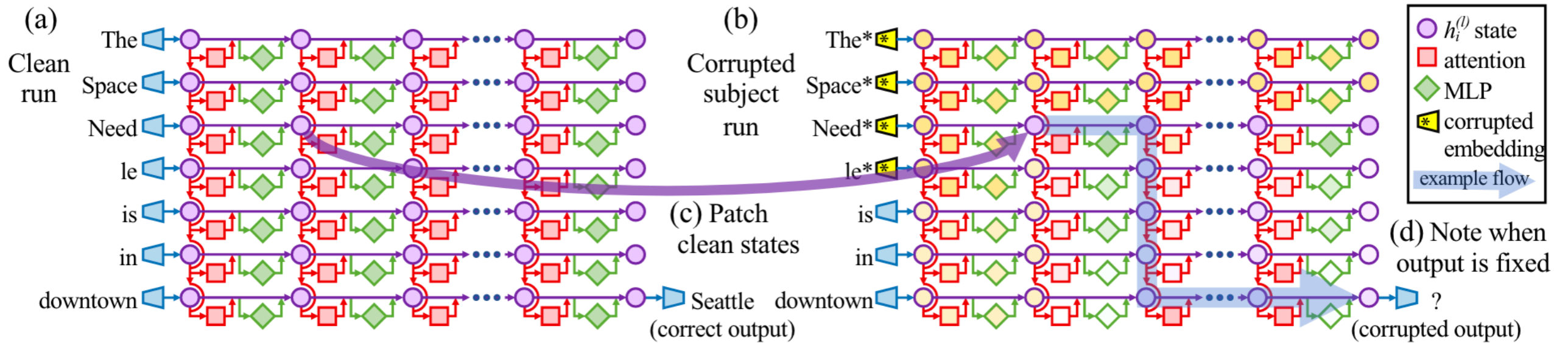(e) Impact of restoring state after corrupted input — early site, late site — single restored layer within GPT-2-XL — p(Seattle)

(f) Impact of restoring MLP after corrupted input — early site — center of interval of 10 restored MLP layers — p(Seattle)

(g) Impact of restoring Attn after corrupted input — late site — center of interval of 10 restored Attn layers — p(Seattle)

[Meng et al., 2022]

# Manipulating key-value associations



[Meng et al., 2022]

# Just ask the model??

| Role | Text |
|------|------|
| User | Can you tell me whether the movie review is positive or negative? <br> The gorgeously elaborate continuation of "The Lord of the Rings" trilogy is so huge that a column of words can not adequately describe co-writer/director Peter Jackson 's expanded vision of J.R.R. Tolkien 's Middle-earth . |
| System | Based on the given text, it appears to be a positive movie review. The phrases "gorgeously elaborate," "expanded vision," and "Middle-earth" suggest admiration and appreciation for the film. |



[Huang et al., 2023]

# Inner workings: overcomplete representations



[Bricken et al., 2023]

# Inner workings: overcomplete representations

#451

**AUTOINTERP. (SCORE = 0.954)**

The neuron fires on the word "so" when used to indicate causation or connect clauses, rather than intensification.

**NEURON ALIGNMENT**

| Neuron | Value | % of $L_1$ |
|---|---|---|
| 454 | +0.27 | 1.8% |
| 305 | +0.22 | 1.5% |
| 400 | +0.20 | 1.4% |

**CORRELATED NEURONS**

| Neuron | Pearson Corr. | Cosine Sim. |
|---|---|---|
| #454 | +0.20 | +0.20 |
| #187 | +0.13 | +0.14 |
| #404 | +0.11 | +0.11 |

**CORRELATED B FEATURES**

| Feature | Pearson Corr. | Cosine Sim. |
|---|---|---|
| #68 | +0.83 | +0.83 |
| #5 | +0.01 | +0.01 |
| #57 | +0.01 | +0.01 |

**ACTIVATIONS (DENSITY = 0.1493%)**

**NEGATIVE LOGITS**

| | |
|---|---|
| eries | −0.38 |
| holder | −0.33 |
| igens | −0.31 |
| NING | −0.30 |
| quirer | −0.29 |
| Aires | −0.29 |
| ción | −0.29 |
| aternity | −0.29 |
| thouse | −0.28 |
| quisition | −0.28 |

**POSITIVE LOGITS**

| | |
|---|---|
| fter | +0.42 |
| aking | +0.38 |
| jour | +0.38 |
| othed | +0.37 |
| forth | +0.37 |
| othes | +0.36 |
| far | +0.36 |
| much | +0.35 |
| aring | +0.35 |
| apy | +0.35 |

**TOP ACTIVATIONS**
**TRAIN TOKEN MAX ACT = 7.651**

same strange boat as **so** many other schools that

a digital model — **so** it may not work

. It's **so** amusing. A family

see humans as just **so** many dollars to be

. It's **so** simple to do and

have at it — **so** he rolled it into

must be protected —– **so** they are bringing in

of the island — **so** we decided to make

Thomasina's **so** naughty↵As

and tremble just **so** in the warmth of

↵It's **so** easy to take things

. They're **so** cute and friendly,

, it's **so** inconceivable that

. Its **so** easy, as Karl

, it's **so** cool now you can

of raw materials— **so** many tons of steel

can be translated as **so** what?↵

and intense story — **so** whether you want a

a cold liquid — **so** what? To me

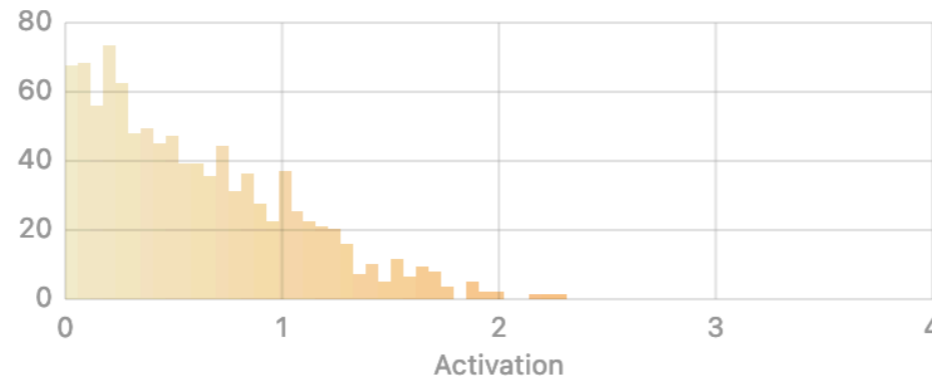, yet with values **so** simple and assured that

# Inner workings: overcomplete representations

#187

**AUTOINTERP. (SCORE = 0.317)** ?

The neuron fires on sentences expressing anticipation, expectation, uncertainty, or doubt.

**ACTIVATIONS (DENSITY = 16.4900%)** ?



**TOP ACTIVATIONS** ?
**TRAIN TOKEN MAX ACT = 3.688**

hands on him before **too** much longer though no

**\xe7\xa7**\x81    \xe8\xa9\xb1
さんが↵        私        の        話

**\xe7\xa7**\x81
?↵「        私        だけ\xe3\x81

as he cannot let **go** of it. On

send for us before **too** much longer."↵

."↵"**Too** bad." He put

nine, was let **go** from his 15–

, then eventually let **go** (see December 16

didn't notice until **too** late that he'd

–Perez as **too** lenient, since

↵VS says "**Too** few arguments...", but

's party to let **go** of it. It

–crush before **too**). Maybe get some

↵↵How to **go** back after accidentally hitting

Same deal as before **too** – predominantly suffixes and

band had to let **go** of their original leader

employees who were let **go** on Monday received a

Q:↵**Too** many fields bad for

not wish to let **go** of the rhetoric of

did not meet until **too** late, and then

**NEGATIVE LOGITS** ?

| | |
|---|---|
| ks | −0.72 |
| consin | −0.62 |
| quit | −0.61 |
| mber | −0.61 |
| can | −0.60 |
| des | −0.59 |
| beit | −0.58 |
| ertain | −0.58 |
| suppl | −0.57 |
| prises | −0.56 |

**POSITIVE LOGITS** ?

| | |
|---|---|
| : | +0.65 |
| :** | +0.58 |
| –*– | +0.54 |
| .: | +0.52 |
| !' | +0.51 |
| .]( | +0.50 |
| ensed | +0.50 |
| time | +0.50 |
| .) | +0.48 |
| githubusercontent | +0.48 |