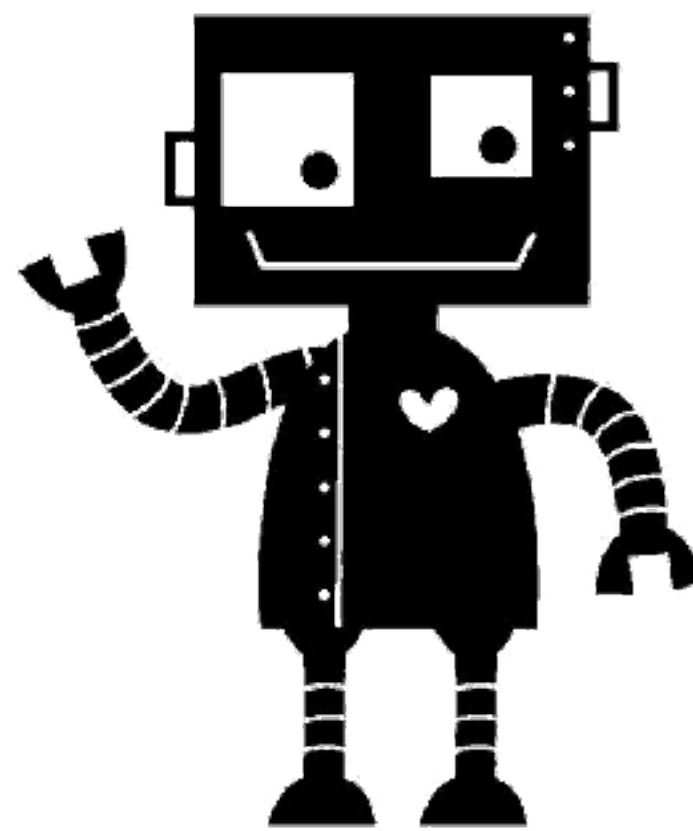


Towards Pragmatic Visual Description Generation

Elisa Kreiss
Department of Communication
UCLA

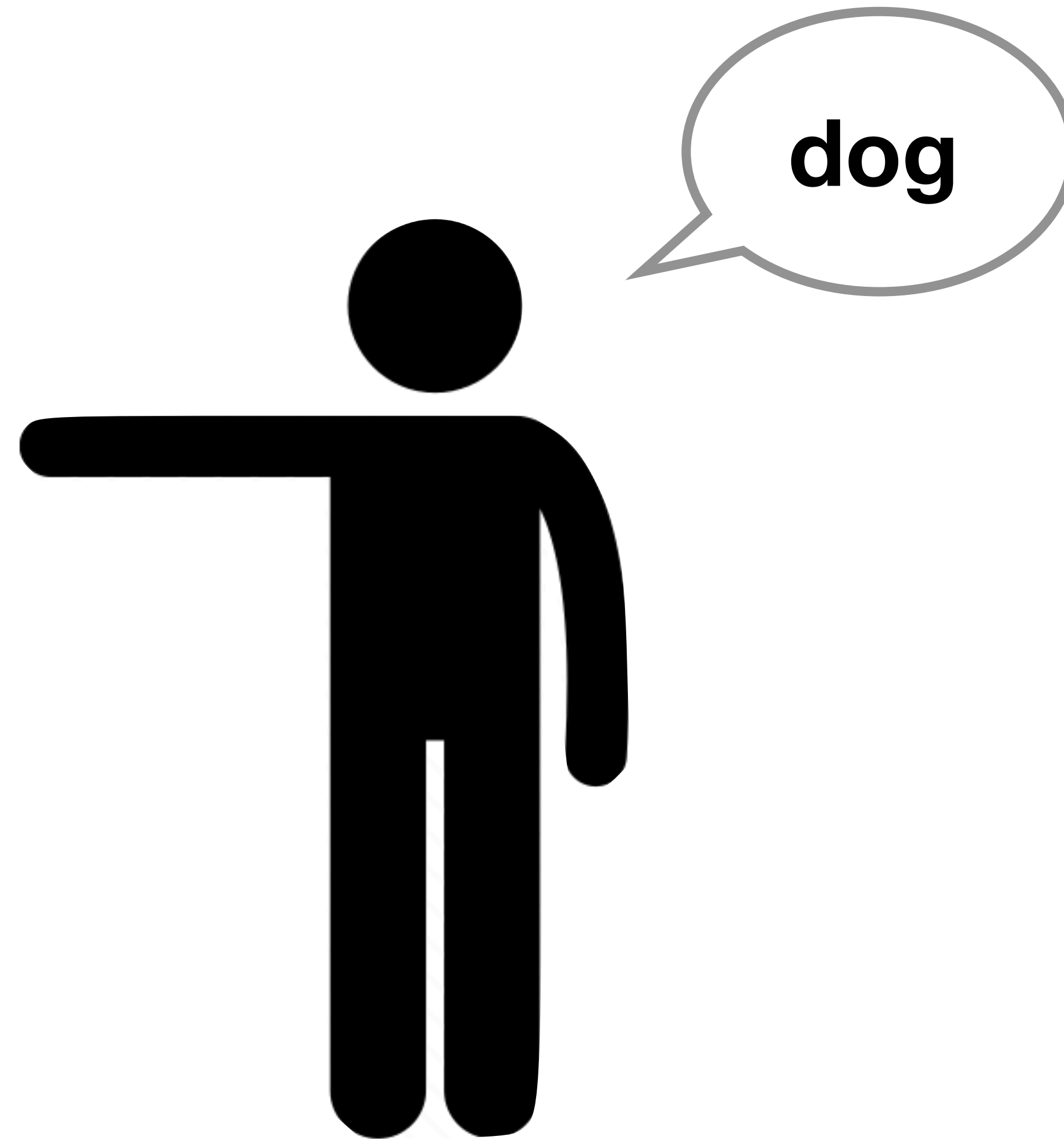


CLASP Seminar
Apr 2 2025



Communicating about a visual world

Communicating about a visual world



Communicating about a visual world

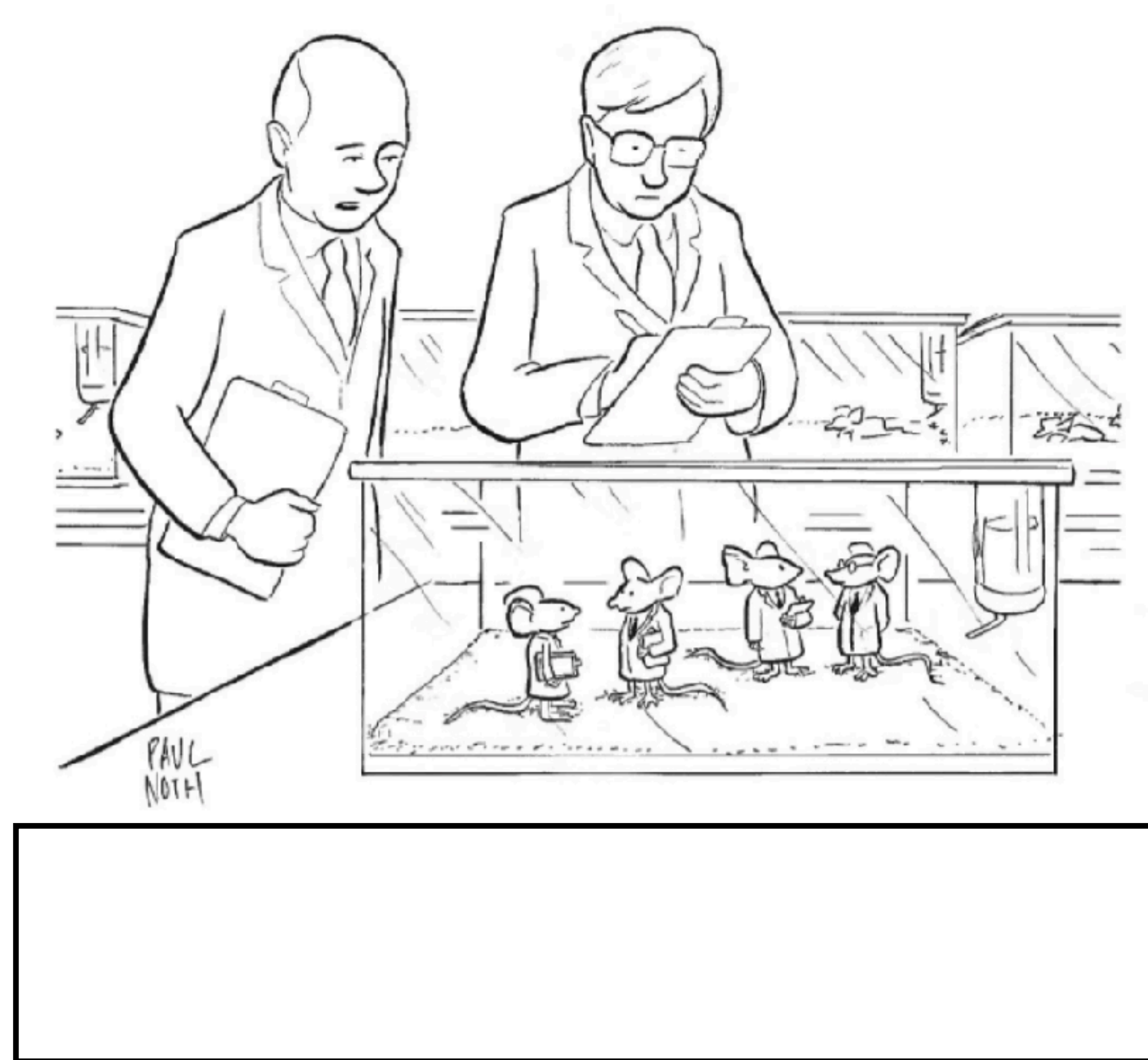


Communicating about a visual world



"De Materie" opera show in New York, starring 100 sheep.

Communicating about a visual world

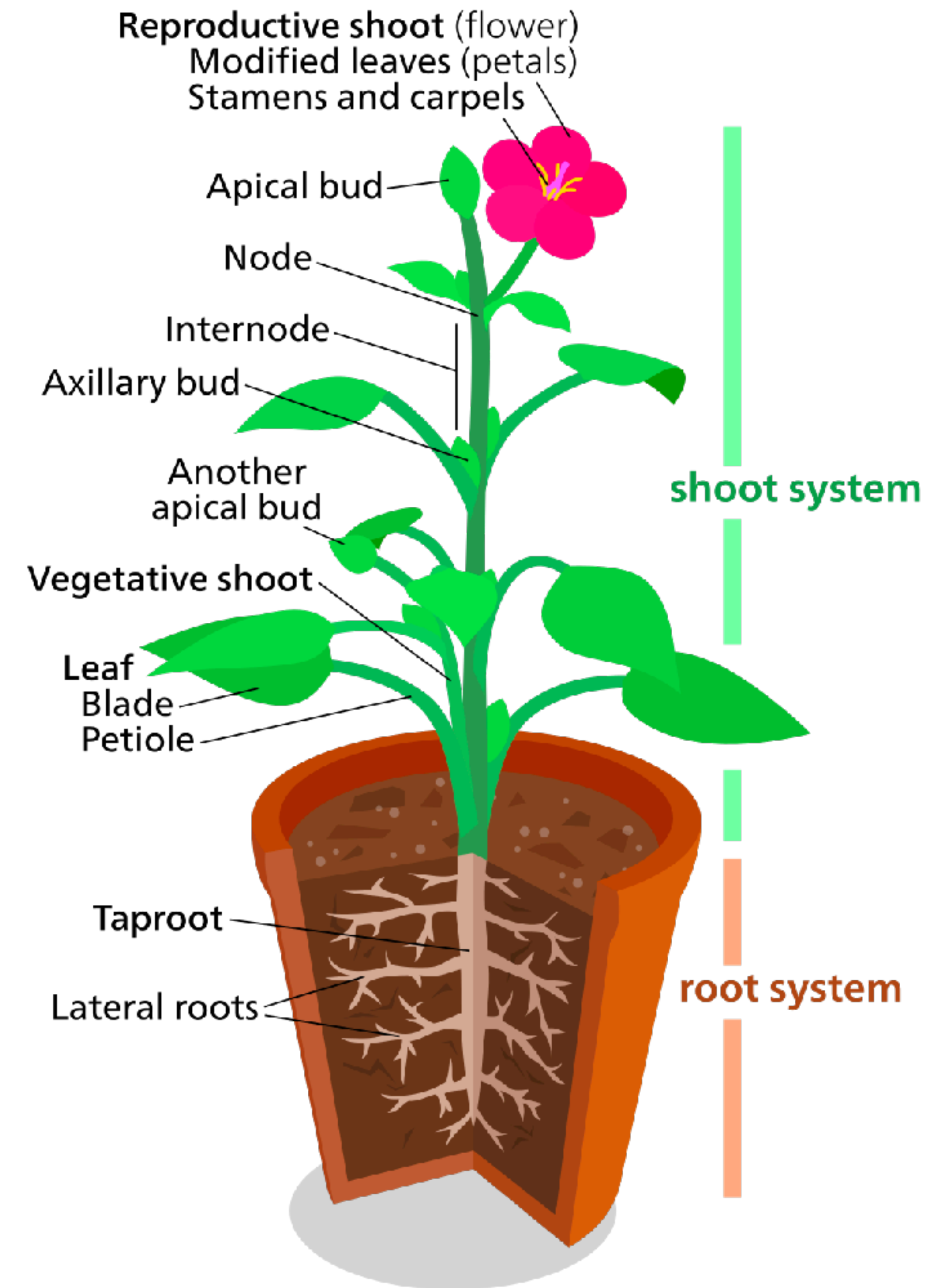


Communicating about a visual world

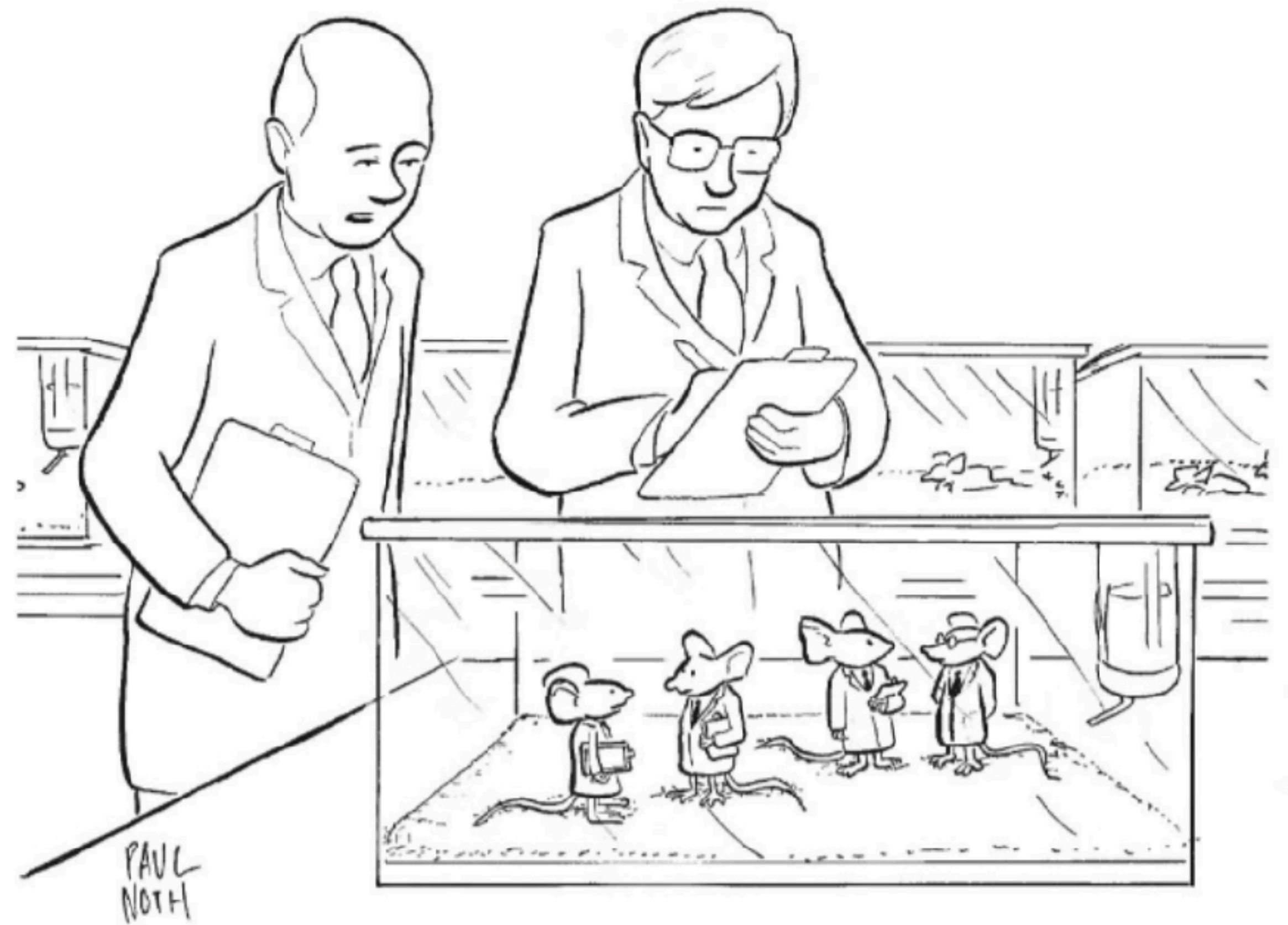


"O.K., let's slowly lower in the grant money."
Todd Bearson, Arlington, Massachusetts.
2009

Education



Entertainment



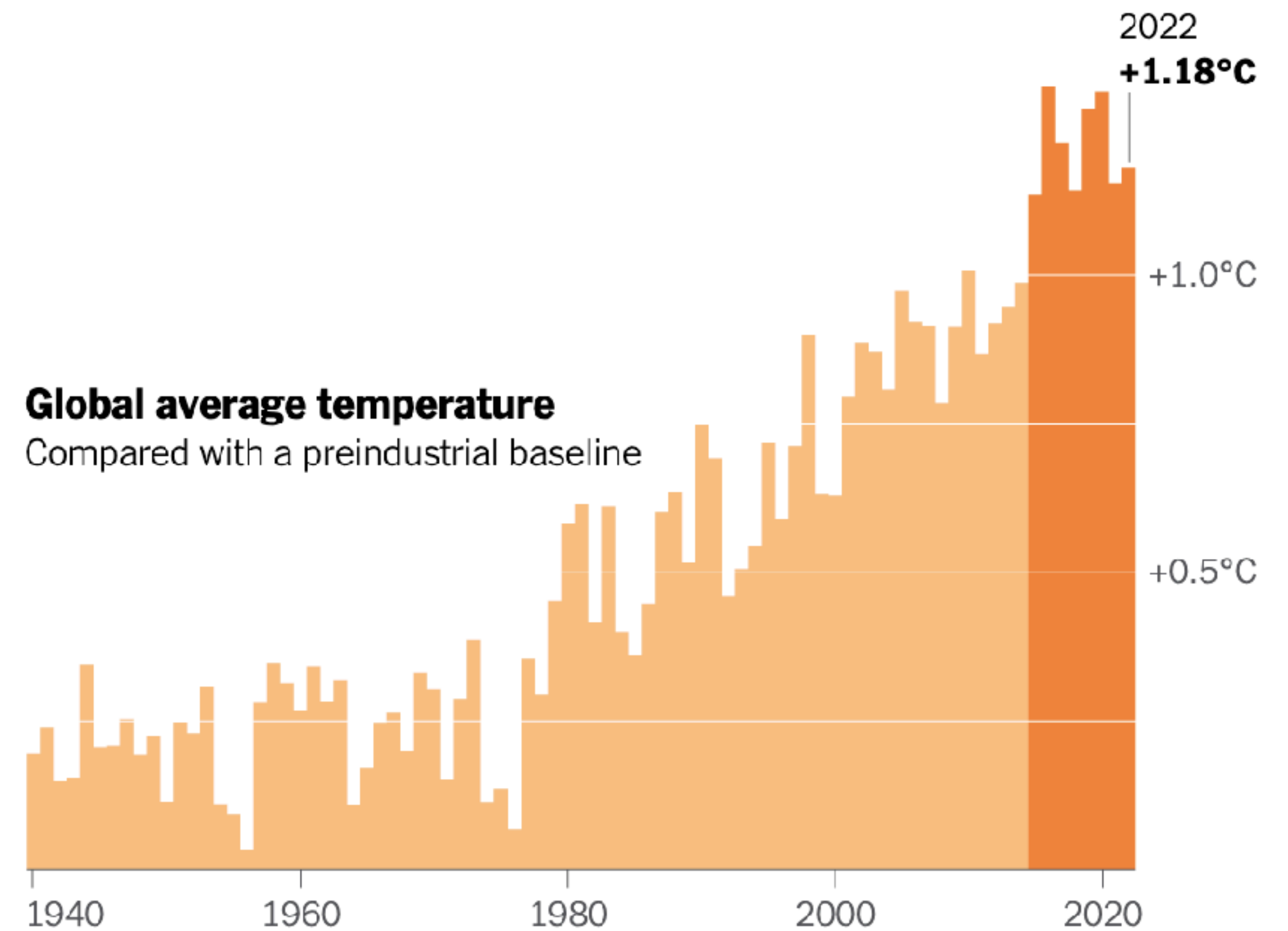
"O.K., let's slowly lower in the grant money."
Todd Bearson, Arlington, Massachusetts.
2009

News

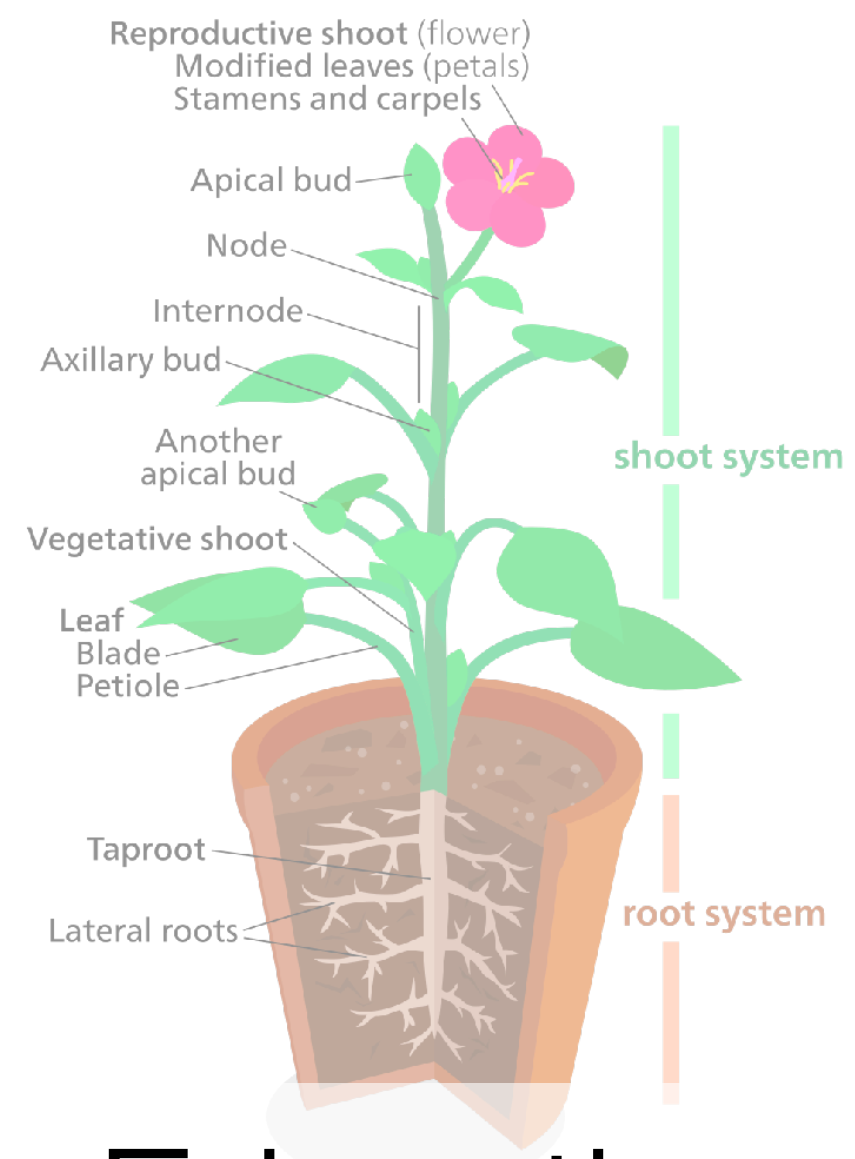
The Last 8 Years Were the Hottest on Record

By [Henry Fountain](#) and [Mira Rojanasakul](#) Jan. 10, 2023

The world remained firmly in warming's grip last year, with extreme summer temperatures in Europe, China and elsewhere contributing to 2022 being the fifth-hottest year on record, European climate researchers said on Tuesday.



Source: Copernicus/ECMWF



Education



Mushroom Embroidery Mushroom Socks Letter...
★★★★★ (1,492)
\$4.14 \$6.28 (50% off)
KYIV

More like this →

Floral Socks, Vintage Socks, Ladies Socks, Wo...
★★★★★ (1,492)
\$4.14 \$6.28 (50% off)
KYIV

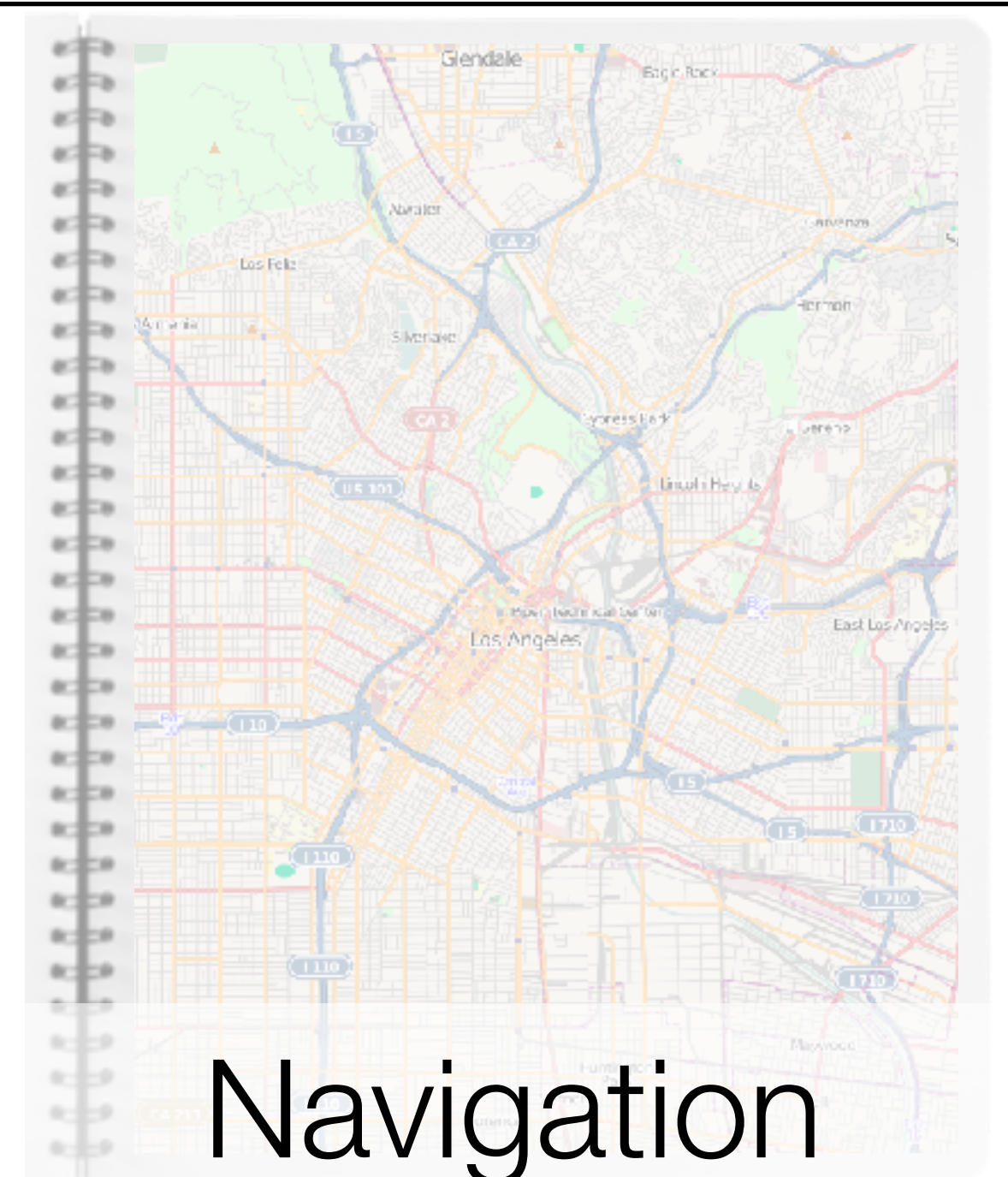
More like this →



Sloth Sock | unisex cozy pastel...
★★★★★ (888)
\$10.50 \$14.69 (28% off)

Shark Crocodile Sha...
★★★★★ (888)
\$10.50 \$14.69 (28% off)

Shopping



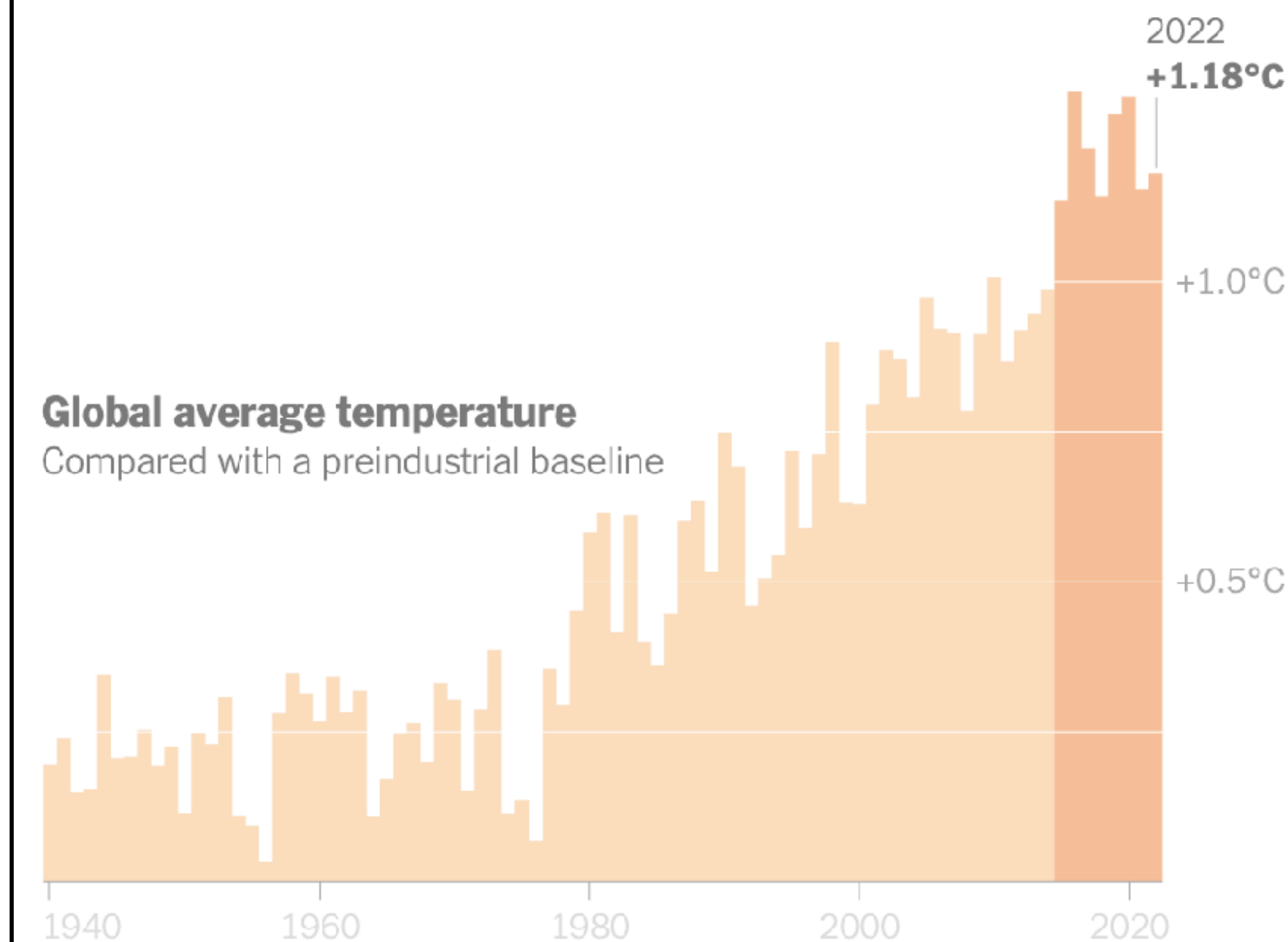
Navigation



"O.K., let's slowly lower in the grant money."
Todd Bearson, Arlington, Massachusetts.
2009

Entertainment

Communication
happens in
complex visual
worlds



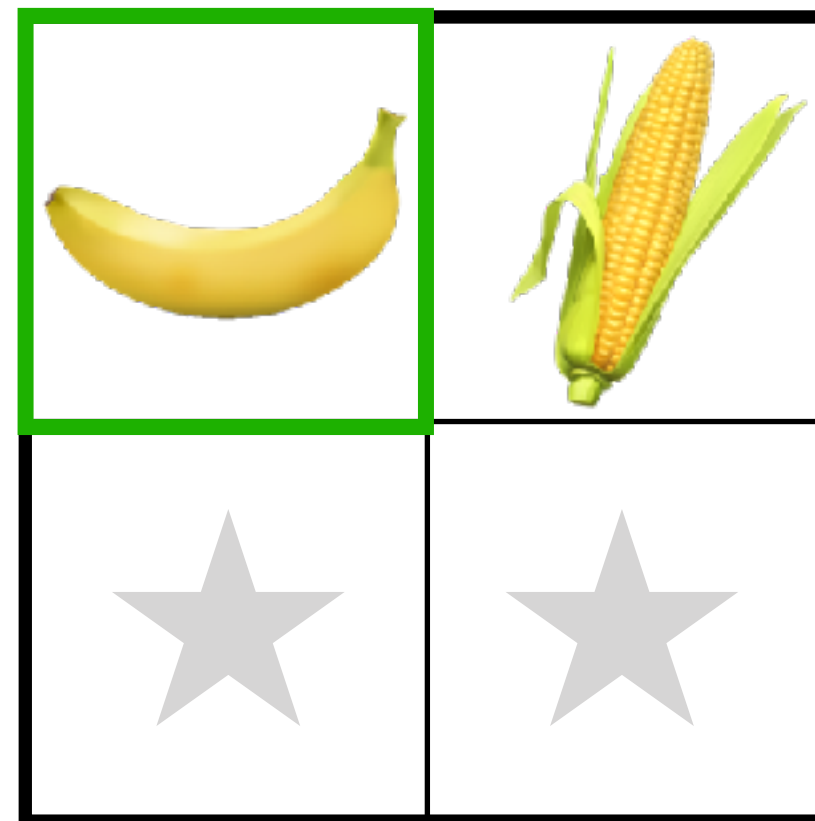
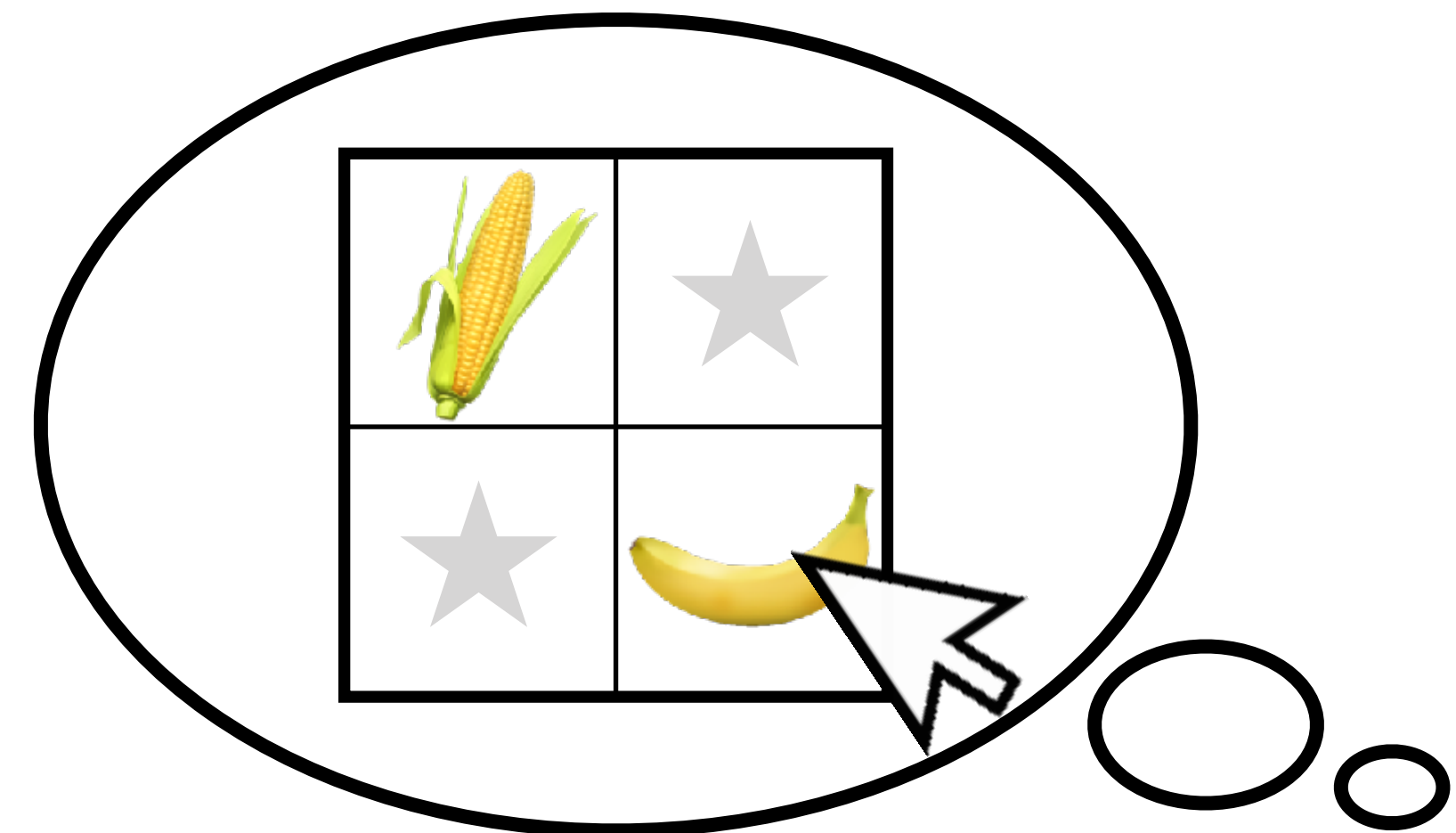
Source: Copernicus/ECMWF

News

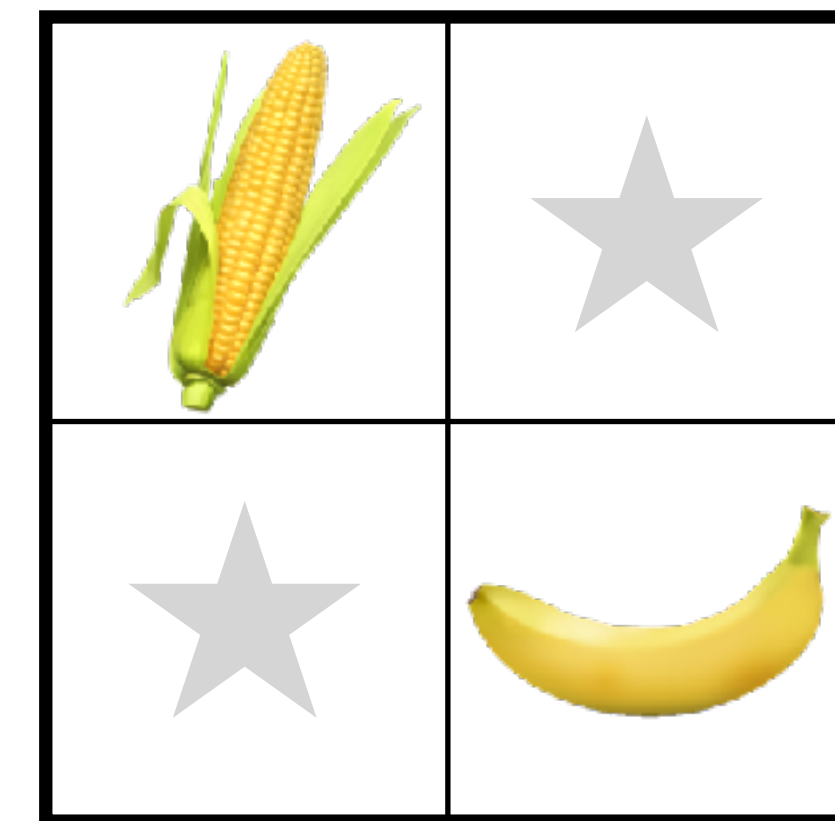


Communication
happens in
complex visual
worlds

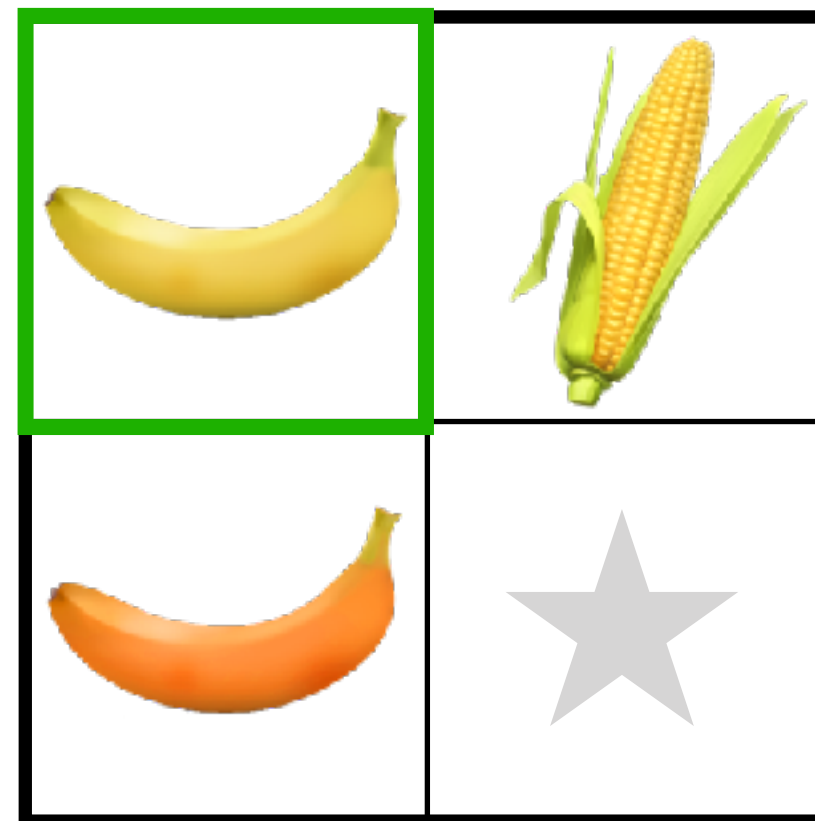
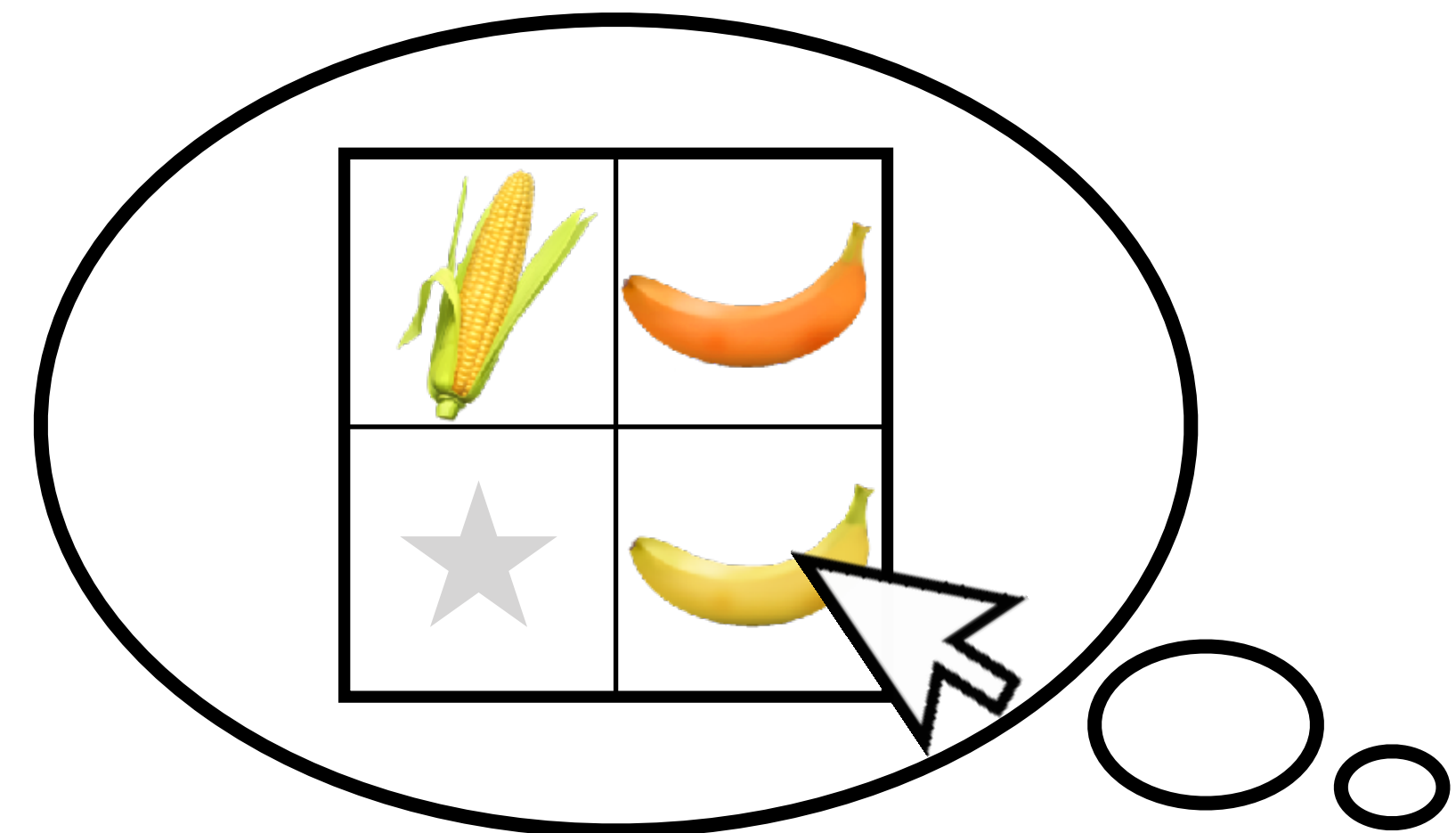
Click on the
banana!



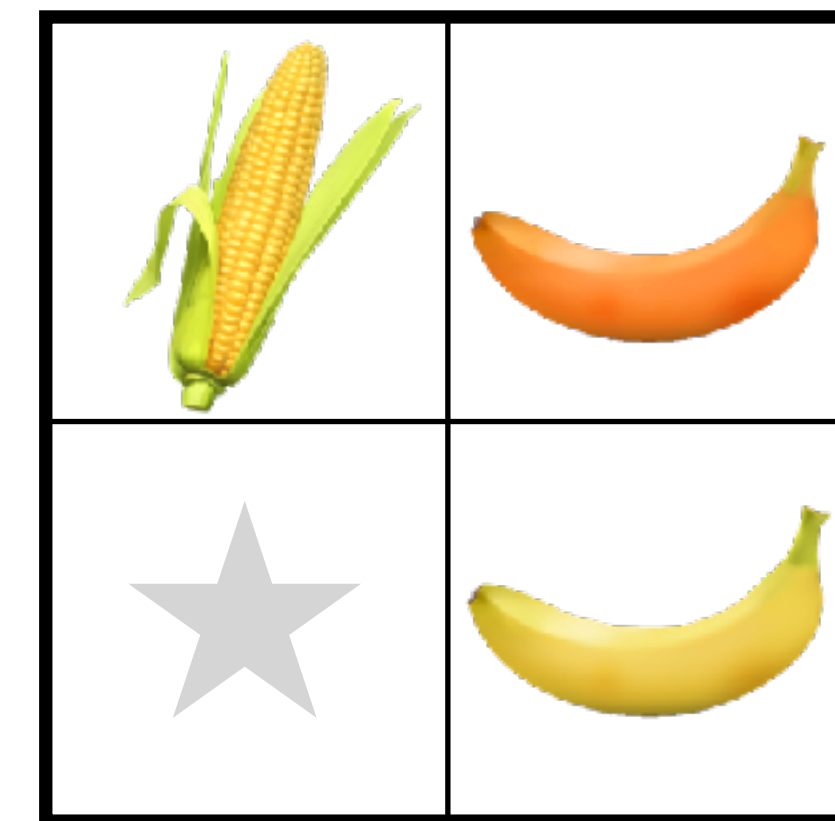
Communication
happens in
complex visual
worlds



Click on the
yellow banana!



Communication
happens in
complex visual
worlds





Communication
happens in
complex visual
worlds

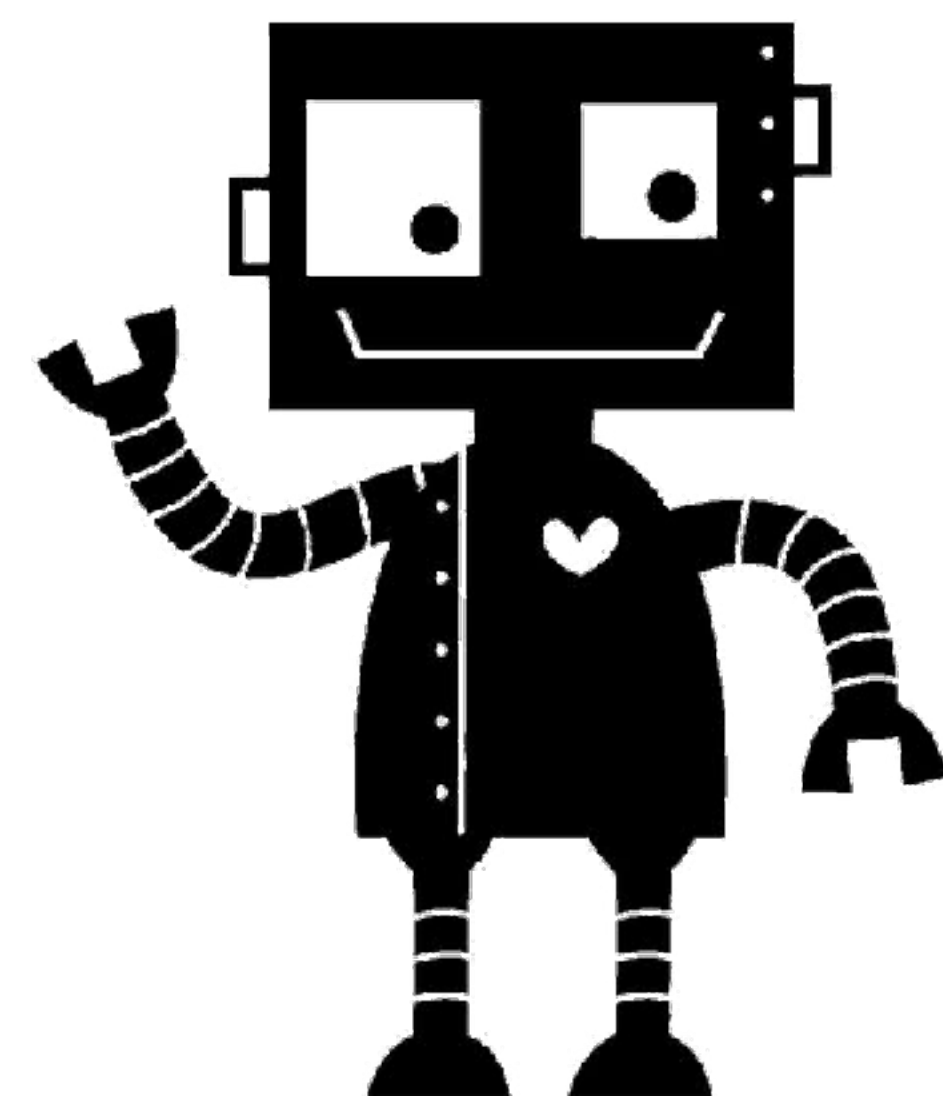


CAT

or

DOG

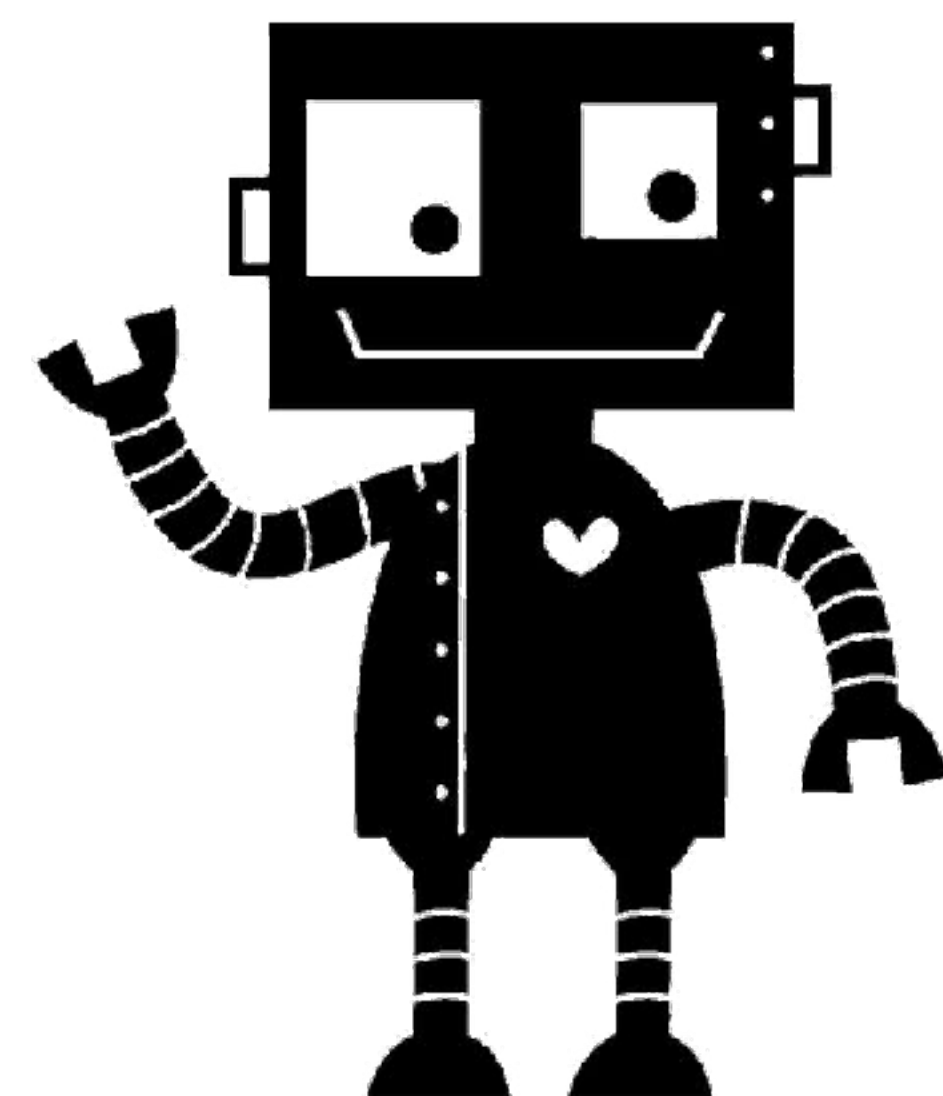
Communication
happens in
complex visual
worlds





"a dog"

Communication
happens in
complex visual
worlds

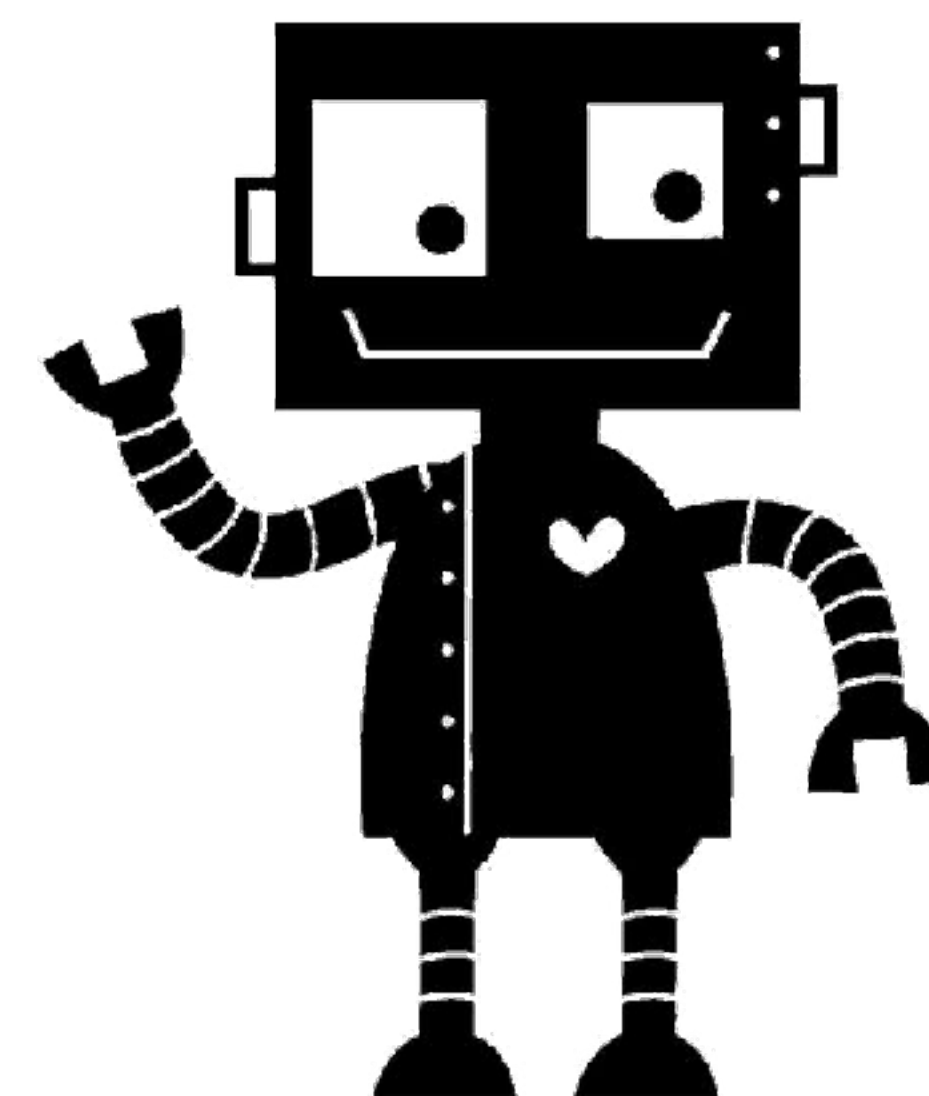


inspired by Karpathy & Fei Fei 2015



**there is a
brown dog
that is lying on
a couch**

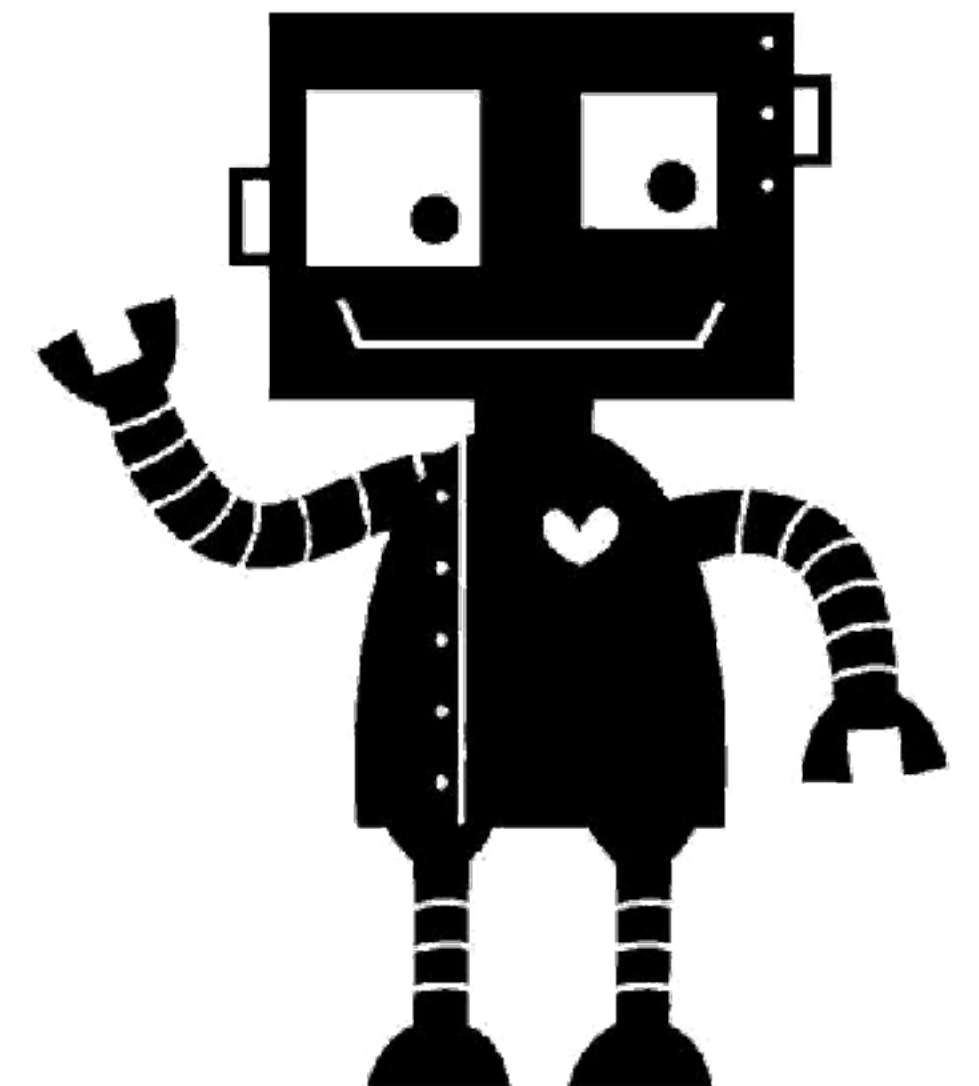
Communication
happens in
complex visual
worlds





**there is a
brown dog
that is lying on
a couch**

Core task





**there is a
brown dog
that is lying on
a couch**

Core task ...

... situated in a
communicative
context

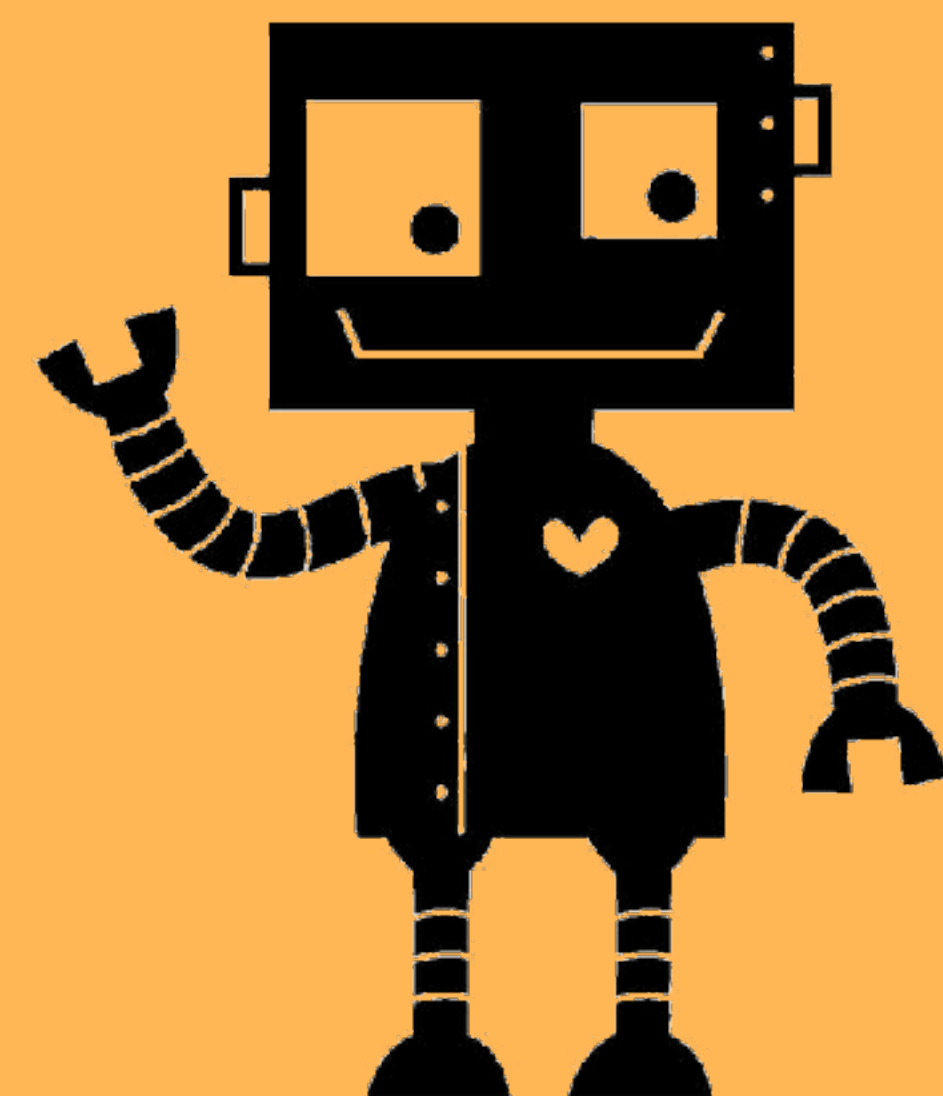
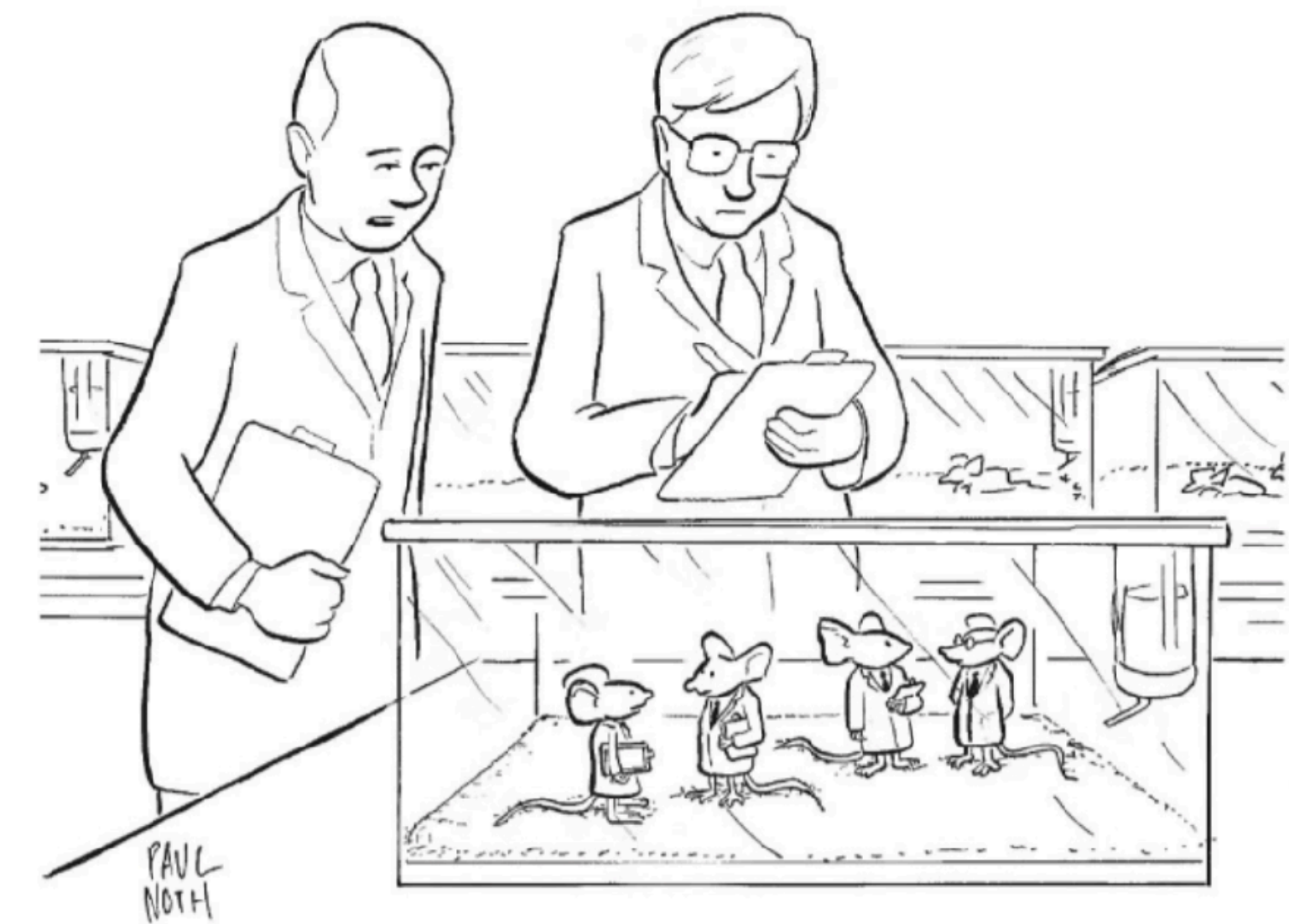


Image Accessibility

An Opportunity and Challenge for AI

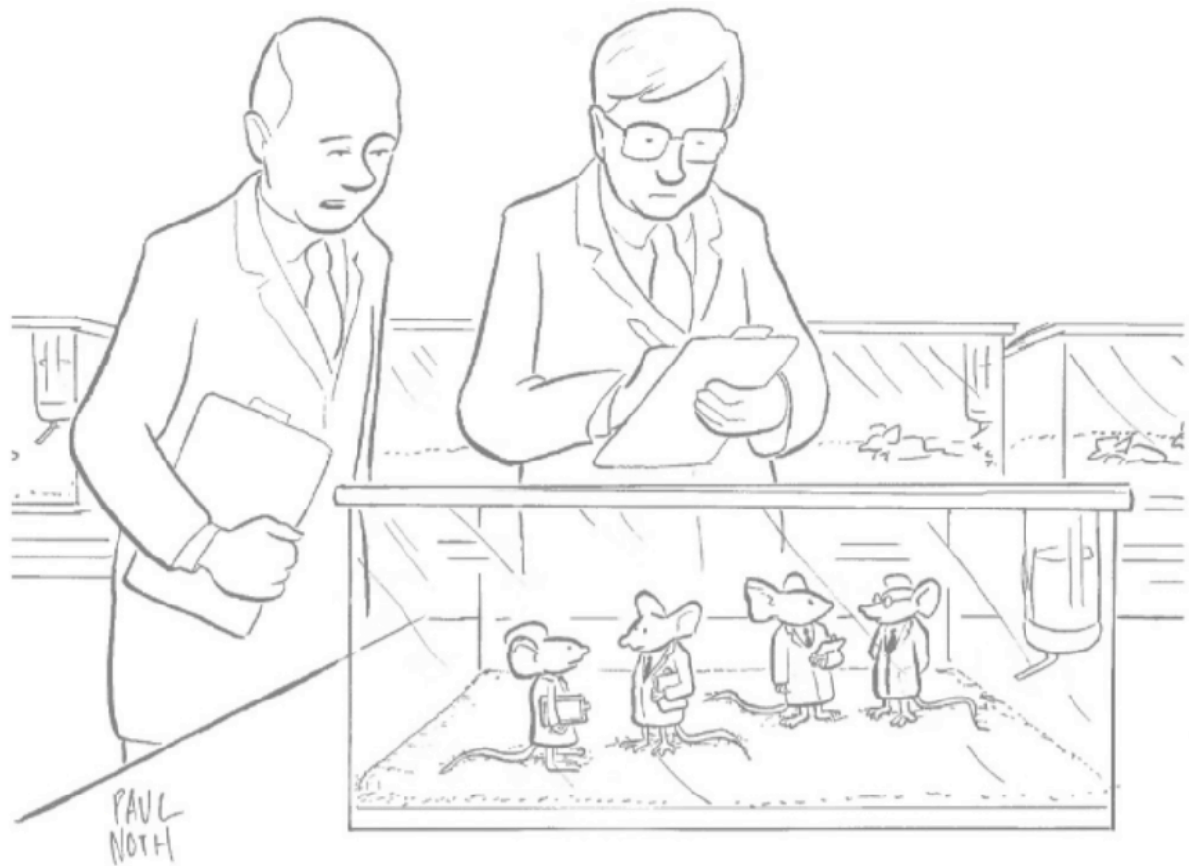
Globally, at least 2.2 billion people have a near or distance vision impairment.
World Health Organization, 2022



"O.K., let's slowly lower in the grant money."
Todd Bearson, Arlington, Massachusetts.
2009

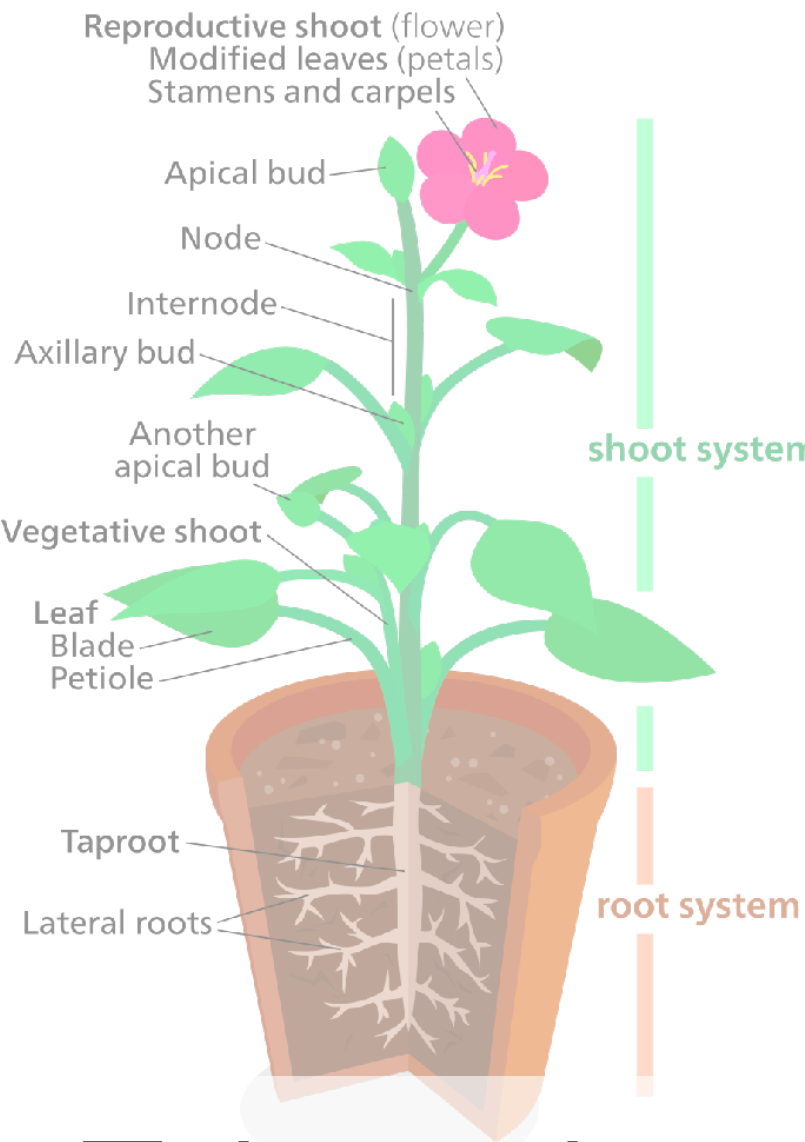
Globally, at least 2.2 billion people have a near or distance vision impairment.
World Health Organization, 2022

The promise of the internet: The leveling playing field for **equal access.**

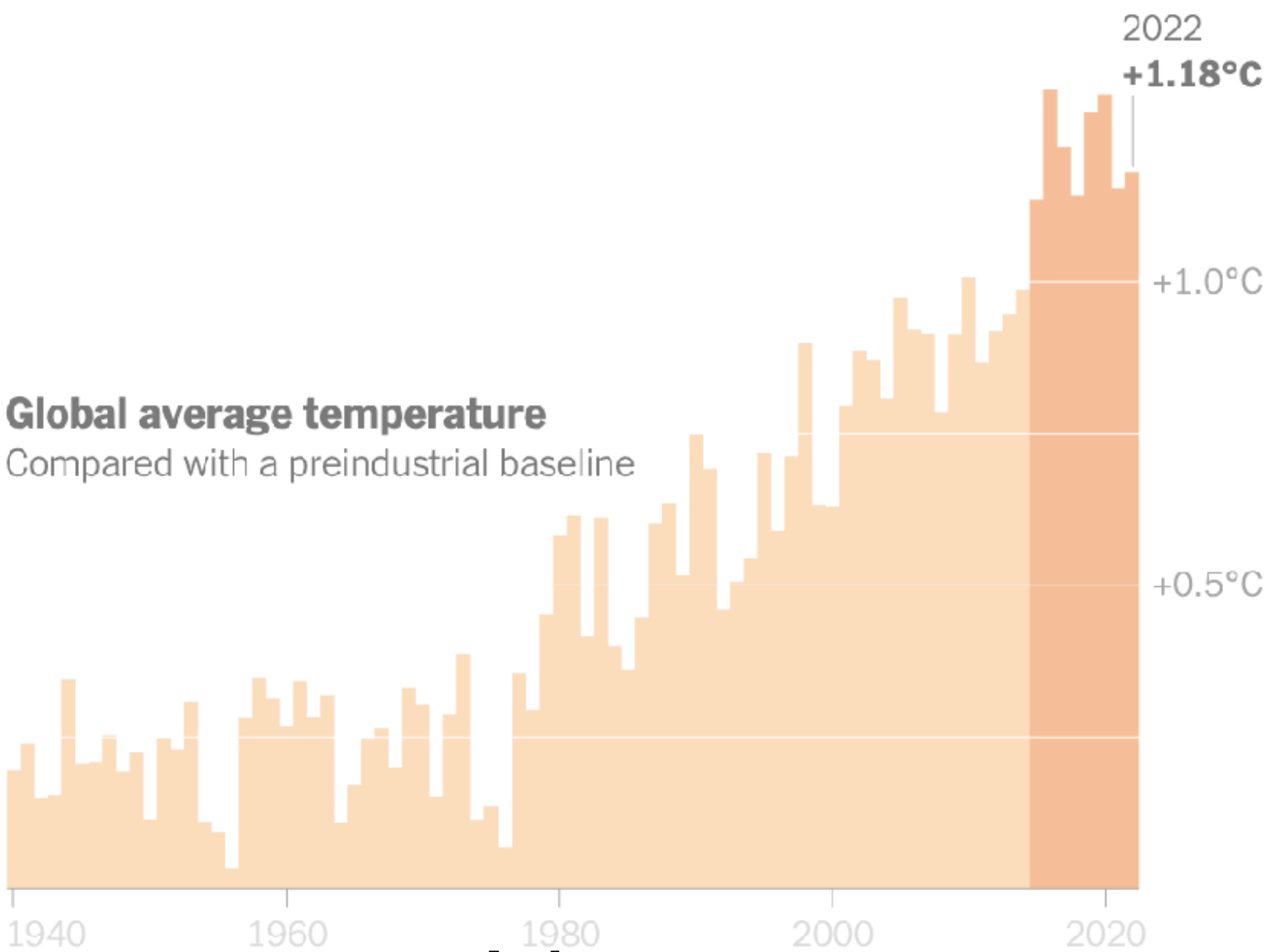


“O.K., let’s slowly lower in the grant money.”
Todd Bearson, Arlington, Massachusetts.
2009

Entertainment



Education



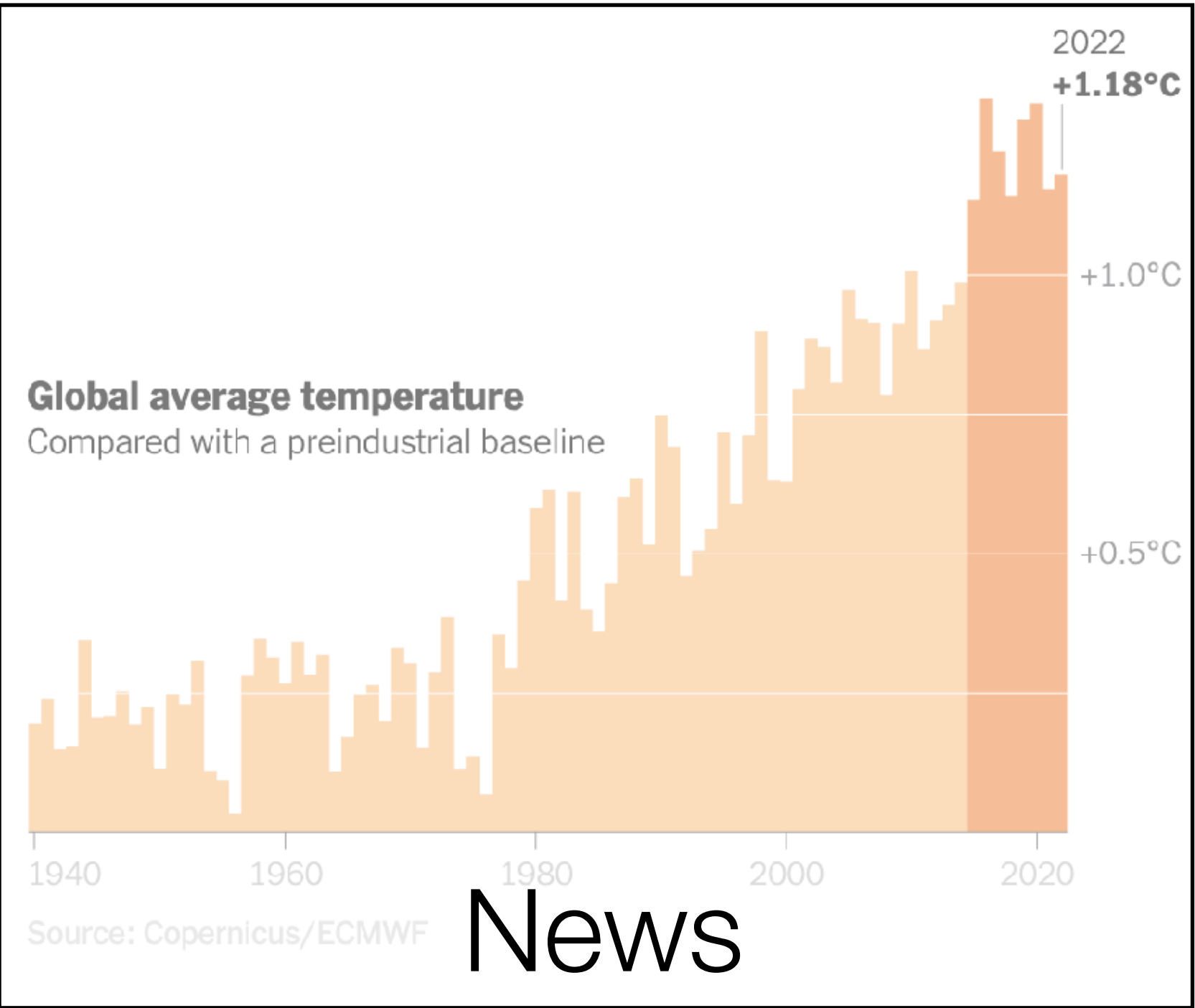
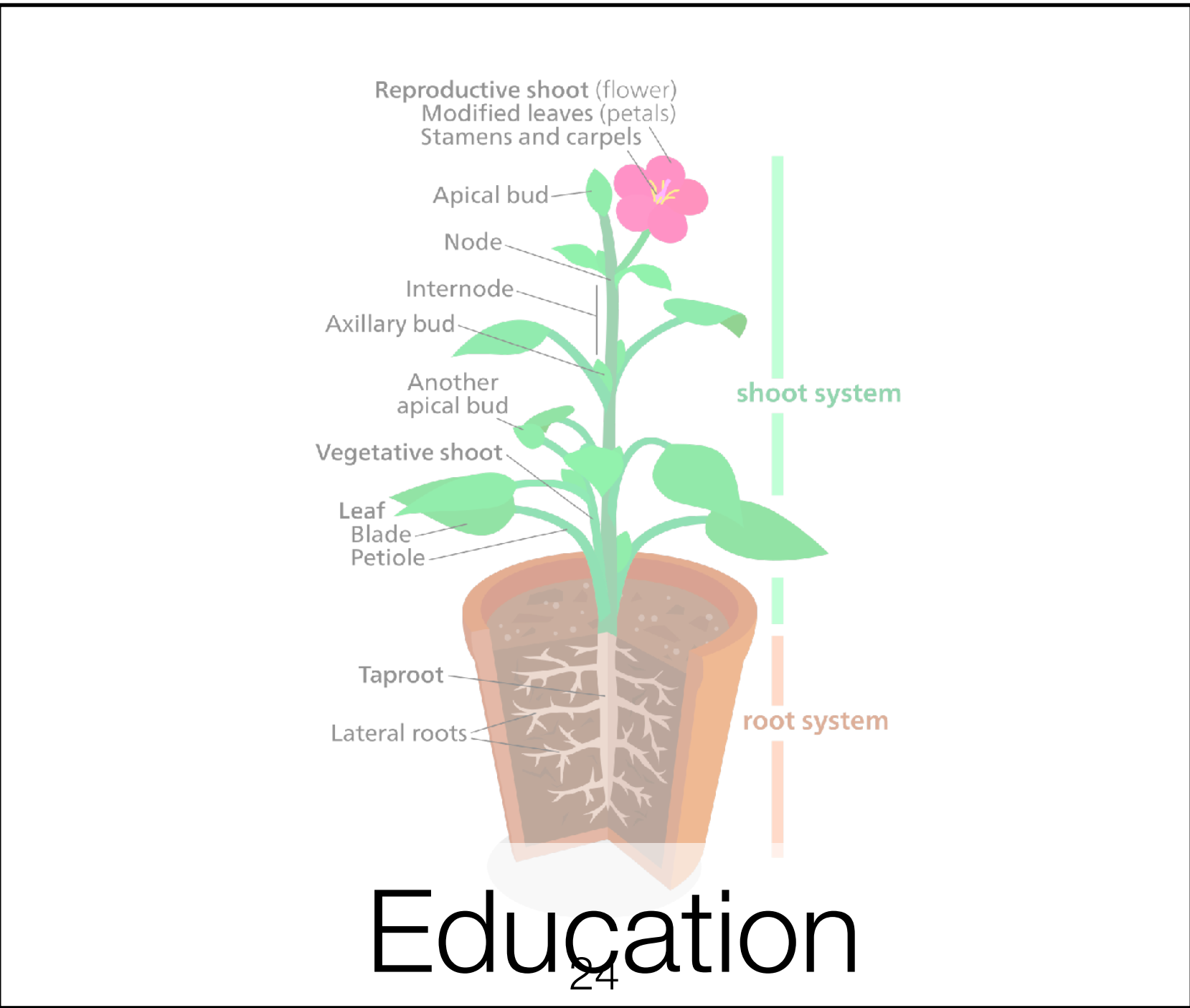
Global average temperature
Compared with a preindustrial baseline

Source: Copernicus/ECMWF

News

Globally, at least 2.2 billion people have a near or distance vision impairment.
World Health Organization, 2022

The promise of the internet: The leveling playing field for **equal access**.
The reality of the internet: Those who can't access, are **disadvantaged**.



Accessing the Web When You Can't See

Screen reader: *software program that allows blind or visually impaired users to read the text that is displayed on the computer screen with a speech synthesizer or braille display* (American Foundation of the Blind, 2022)



Elisa Kreiss @ElisaKreiss

I'm so excited for the talk today!



I'm so excited for the talk today!

Accessing the Web When You Can't See

Screen reader: *software program that allows blind or visually impaired users to read the text that is displayed on the computer screen with a speech synthesizer or braille display* (American Foundation of the Blind, 2022)



Elisa Kreiss @ElisaKreiss



Accessing the Web When You Can't See

Screen reader: *software program that allows blind or visually impaired users to read the text that is displayed on the computer screen with a speech synthesizer or braille display* (American Foundation of the Blind, 2022)



Accessing the Web When You Can't See

Screen reader: *software program that allows blind or visually impaired users to read the text that is displayed on the computer screen with a speech synthesizer or braille display* (American Foundation of the Blind, 2022)



Elisa Kreiss @ElisaKreiss

Alt description:

Cute, small dog sitting on a sidewalk, looking up with big eyes. Ears are propped up.



Cute, small dog sitting on a sidewalk, looking up with big eyes. Ears are propped up.

Accessing the Web When You Can't See

Screen reader: *software program that allows blind or visually impaired users to read the text that is displayed on the computer screen with a speech synthesizer or braille display* (American Foundation of the Blind, 2022)



Elisa Kreiss @ElisaKreiss

Alt description:

Cute, small dog sitting on a sidewalk, looking up with big eyes. Ears are propped up.



Cute, small dog sitting on a sidewalk, looking up with big eyes. Ears are propped up.



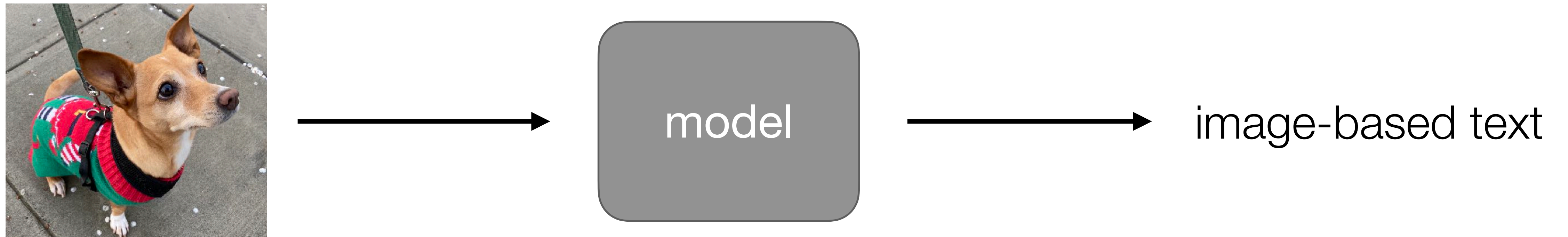
very rare

e.g., only 0.1% of Twitter images have alt text

Opportunity for AI to help!

Visual Accessibility: **An Opportunity for AI**

Vision-Language models translate images into text.



Visual Accessibility: **An Opportunity for AI**

Image-based text generation

Vision-Language models translate images into text.



model



image-based text

Visual Accessibility: **An Opportunity for AI**

Image-based text generation



A man in a suit gestures while speaking, with the U.S. flag in the background.

(GPT-4V, OpenAI 2023)

Visual Accessibility: **An Opportunity for AI**

There are many texts that can go with every single image.

Barack Obama giving his last state of the union address.

A hand with an index finger that points up, seemingly attached to a person who has their mouth open. Teeth are visible.

...

A microphone facing away from the viewer. A flag is in the background. A man stands between the microphone and the flag.



A man in a suit gestures while speaking, with the U.S. flag in the background.

(GPT-4V, OpenAI 2023)

Visual Accessibility: **An Opportunity for AI**

Challenge: What makes an accessibility description useful?

Barack Obama giving his last state of the union address.

A hand with an index finger that points up, seemingly attached to a person who has their mouth open. Teeth are visible.

...

A microphone facing away from the viewer. A flag is in the background. A man stands between the microphone and the flag.



A man in a suit gestures while speaking, with the U.S. flag in the background.

(GPT-4V, OpenAI 2023)

Visual Accessibility: **An Opportunity for AI**

Consequence: Even sophisticated model can miss accessibility goals.

[**See also:** MacLeod et al., 2017; Bennett et al., 2021; Herskovitz et al., 2023]

Example

Be My AI: OpenAI / Be My Eyes
(GPT-4V(ision) System Card, 2023)



A man in a suit gestures while speaking, with the U.S. flag in the background.

(GPT-4V, OpenAI 2023)

Visual Accessibility: **An Opportunity for AI**

Consequence: Even sophisticated model can miss accessibility goals.

[**See also:** MacLeod et al., 2017; Bennett et al., 2021; Herskovitz et al., 2023]

Example

Be My AI: OpenAI / Be My Eyes
(GPT-4V(ision) System Card, 2023)



Barack Obama

?

A man in a suit gestures while speaking, with the U.S. flag in the background.

(GPT-4V, OpenAI 2023)

Reframing Image-Based Text Generation

What **can** we say
about an image?



What **should** we say
about an image?

Reframing Image-Based Text Generation

What **can** we say
about an image?



What **should** we say
about an image?

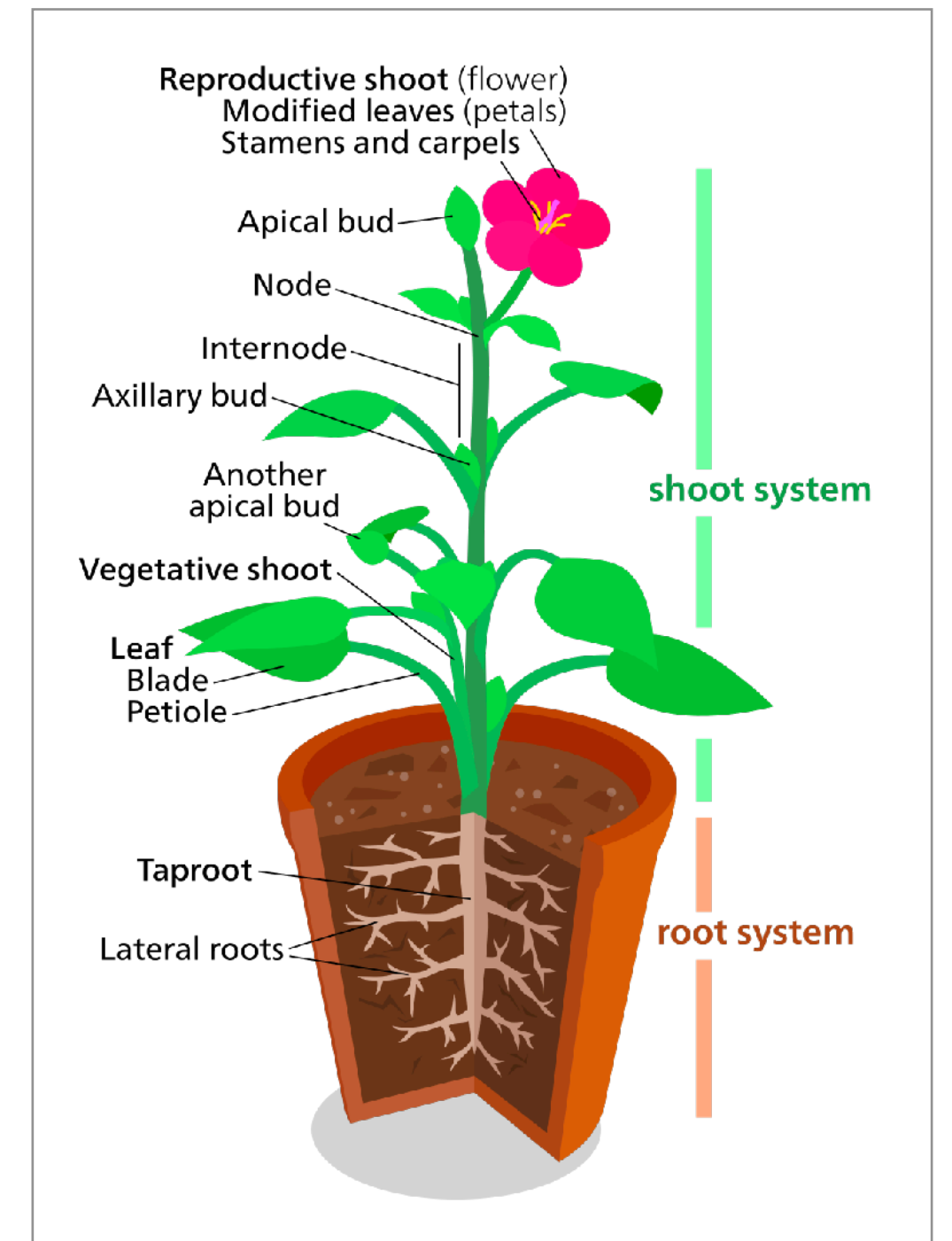
ground-truth description → **pragmatically useful** description

Image-based text generation depends on ...

1 the **image-based text**'s communicative goal.

2 the **image**'s communicative goal.

The Communicative Goal of the Image-Based **Text**





WIKIPEDIA
The Free Encyclopedia

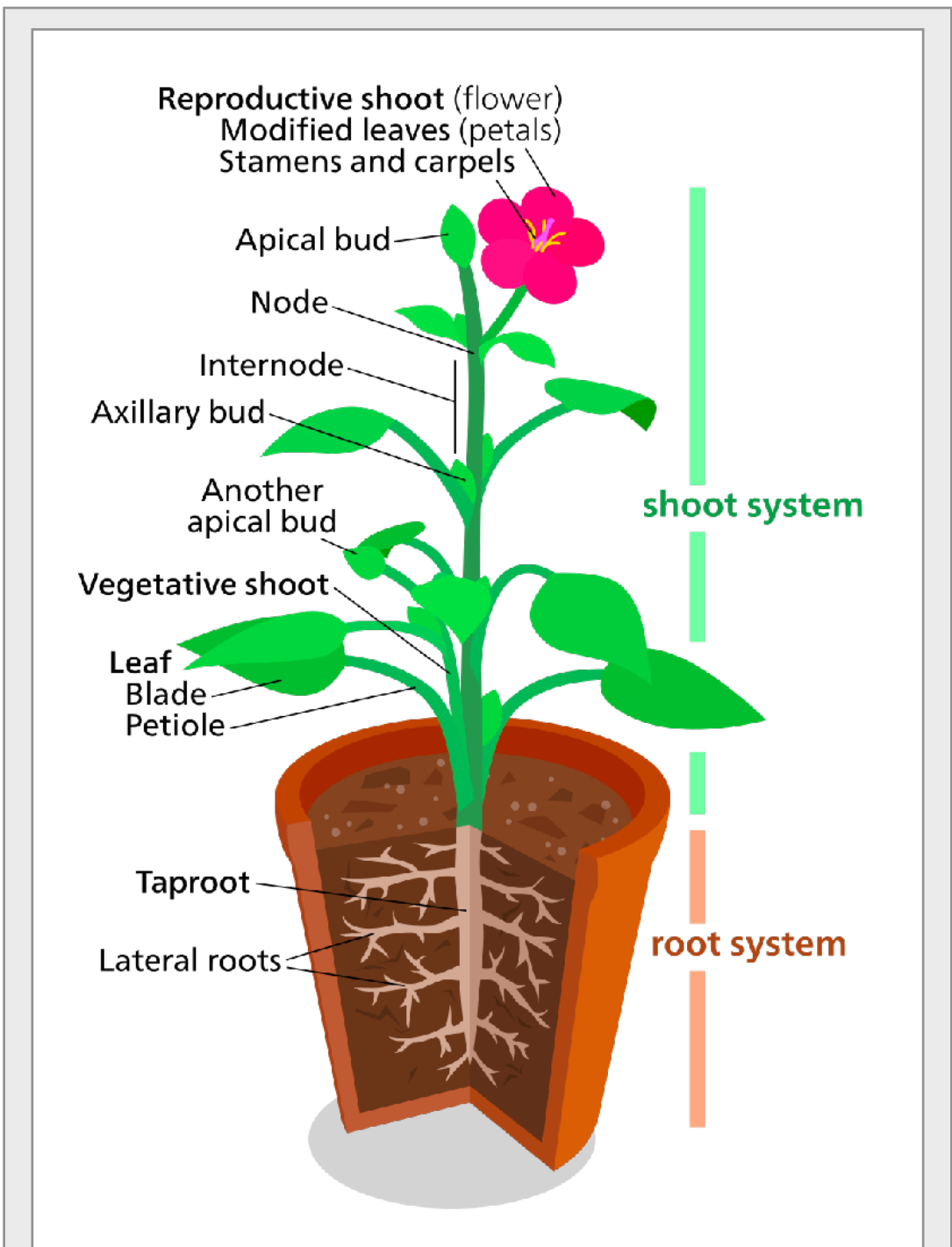
The Communicative Goal of the Image-Based Text

Multimodal Pedagogy

Multimodal pedagogy is an approach to the teaching of writing that implements different modes of communication.^{[1][2]} **Multimodality** refers to the use of visual, aural, linguistic, spatial, and gestural modes in differing pieces of media, each necessary to properly convey the information it presents.^{[3][4]}

The visual mode conveys meaning via images and the visible elements of a text such as typography and color. The aural mode refers to sound in the form of music, sound effects, silence, etc. The linguistic mode includes written and spoken language. The spatial mode focuses on the physical arrangement of elements in a text. The gestural mode refers to physical movements such facial expressions and how these are interpreted. A multimodal text is characterized by the combination of any two or more modes to express meaning.^[5]

Multimodality as a term was coined in the late 20th century,^[6] but its use predates its naming, with it being used as early as **Egyptian hieroglyphs** and classical **rhetoric**.^[7] Compositionists and writing theorists have been exploring how the five modes of communication interact with each other and how multimodality can be used in the teaching of writing since the 20th century.^[8]



Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.



WIKIPEDIA
The Free Encyclopedia

The Communicative Goal of the Image-Based Text

Multimodal Pedagogy

Multimodal pedagogy is an approach to the teaching of writing that implements different modes of communication.^{[1][2]} **Multimodality** refers to the use of multiple modes of communication, such as text, images, and sound, in a single communication.

The visual mode conveys meaning through images, diagrams, and other visual elements.

The aural mode refers to sound in the form of spoken language or audio recordings.

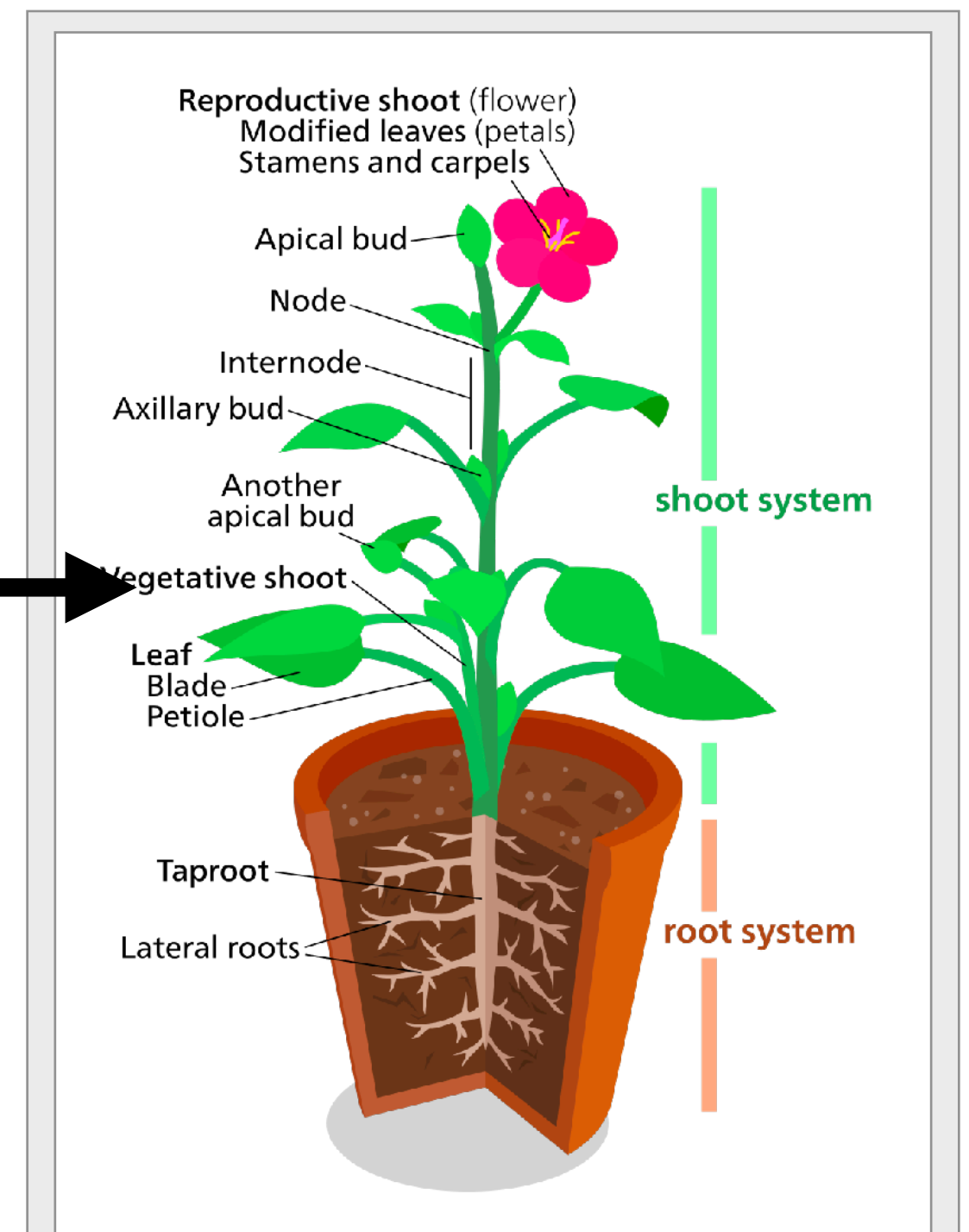
The written mode refers to text in the form of printed or digital documents.

The gestural mode refers to physical movement or gestures that convey meaning.

Text is characterized by the combination of these modes, creating a multimodal communication.

Multimodality as a term was coined in the 1990s, but has been used as early as **Egyptian hieroglyphs** and classical **rhetoric**.^[7] Compositionists and writing theorists have been exploring how the five modes of communication interact with each other and how multimodality can be used in the teaching of writing since the 20th century.^[8]

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.



Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.



WIKIPEDIA
The Free Encyclopedia

The Communicative Goal of the Image-Based Text

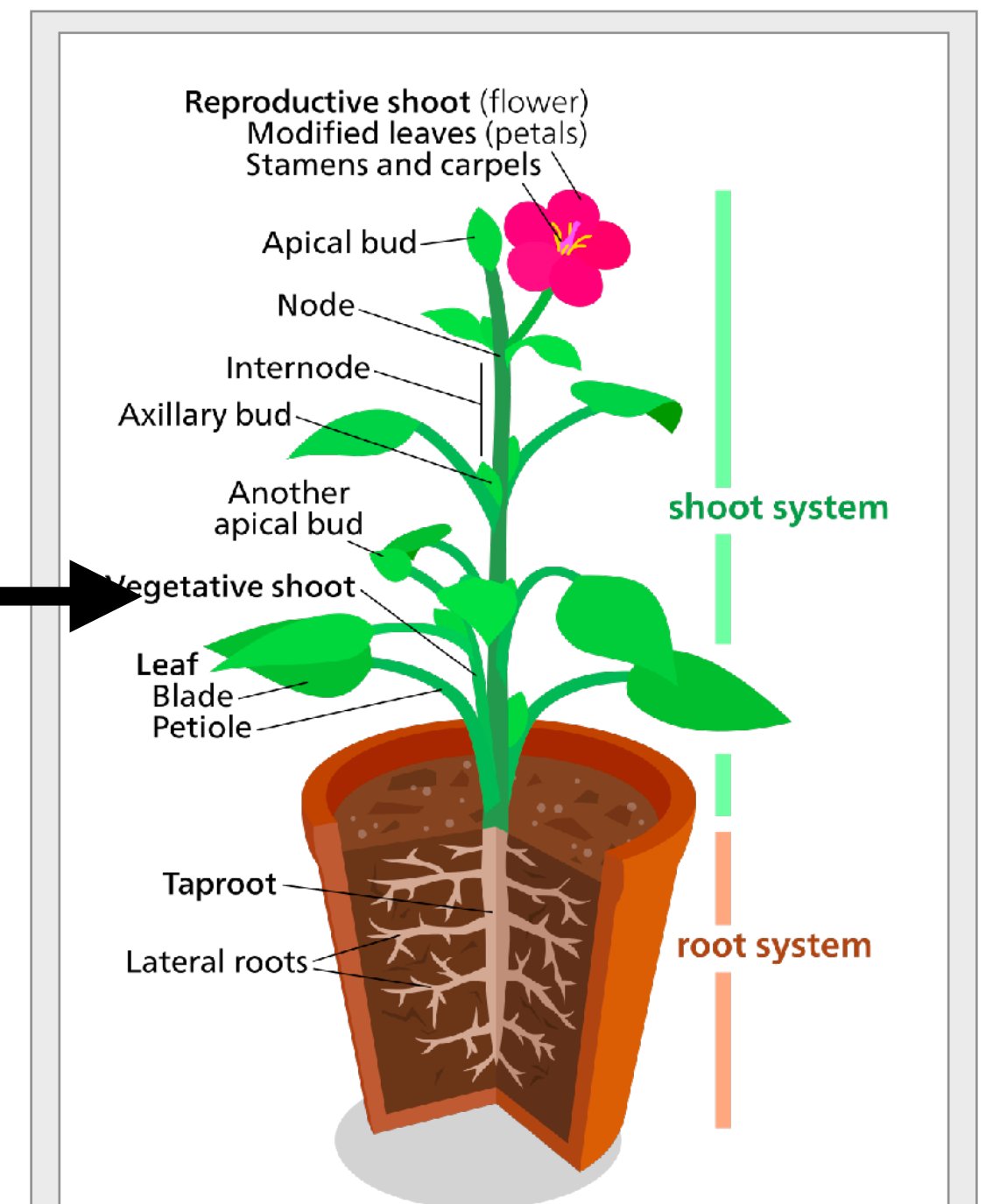
Multimodal Pedagogy

Multimodal pedagogy is an approach to the teaching of writing that implements different modes of communication.^{[1][2]} **Multimodality** refers to the use of different modes of communication, such as text, image, sound, and gesture, in a single communicative act.

The visual mode conveys meaning through images, diagrams, and other visual elements. The aural mode refers to sound in written and spoken language. The gestural mode refers to physical movement. Multimodal text is characterized by the combination of these different modes of communication.

Multimodality as a term was coined in the 1990s and has since been used as early as **Egyptian hieroglyphs** and classical **rhetoric**.^[7] Compositionists and writing theorists have been exploring how the five modes of communication interact with each other and how multimodality can be used in the teaching of writing since the 20th century.^[8]

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.



Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.



WIKIPEDIA
The Free Encyclopedia

The Communicative Goal of the Image-Based Text

Multimodal Pedagogy

Multimodal pedagogy is an approach

communication.^{[1][2]} Multimodality refers

differing pieces of media, each needing

The visual mode conveys meaning

The aural mode refers to sound in

written and spoken language. The

gestural mode refers to physical movement

text is characterized by the combination

Multimodality as a term was coined in

early as **Egyptian hieroglyphs** and classical

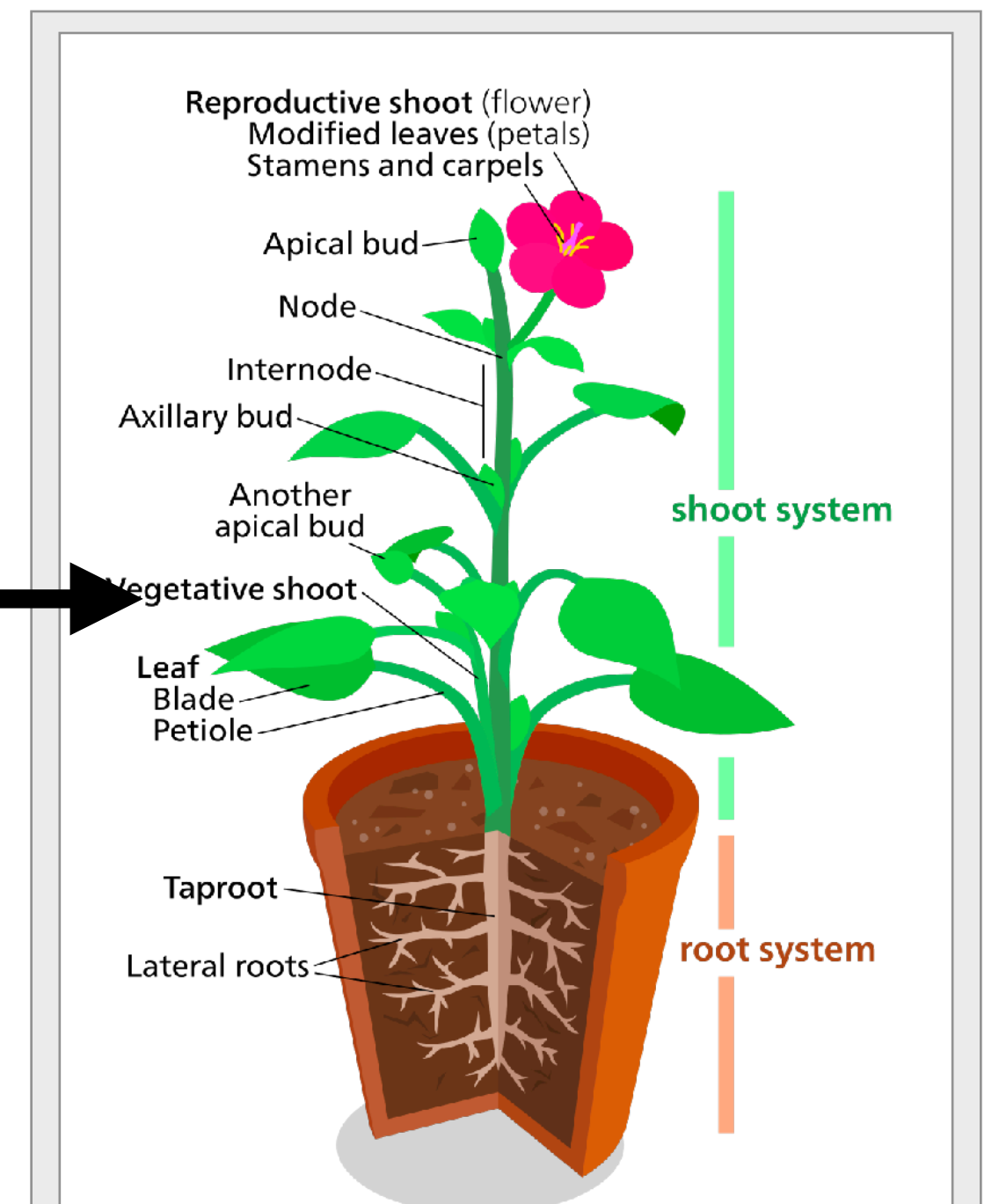
how the five modes of communication interact with each other and how multimodality can be used in the teaching

of writing since the 20th century.^[8]

Description

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil.

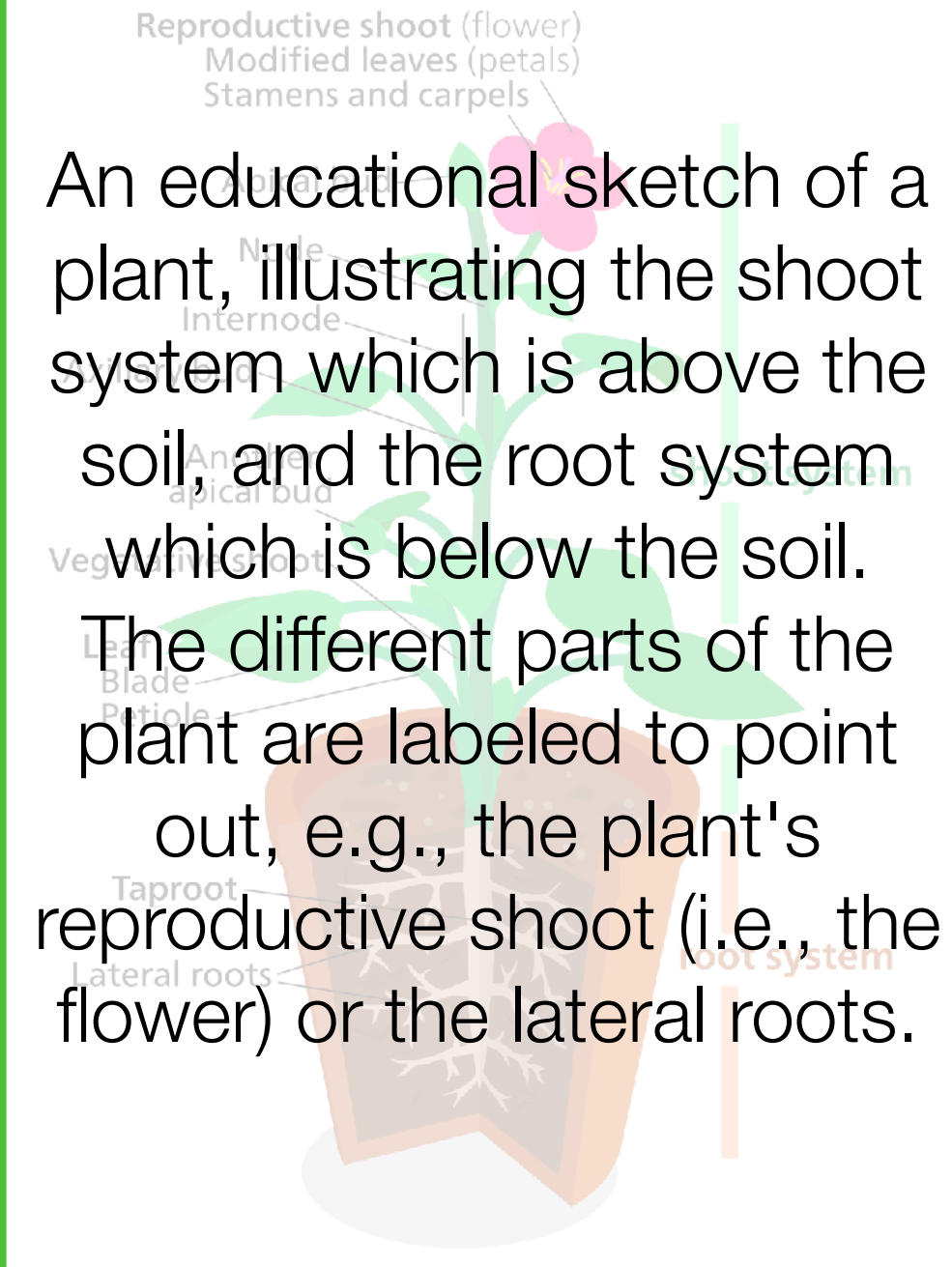
The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.



Caption

Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.

The Communicative Goal of the Image-Based **Text**



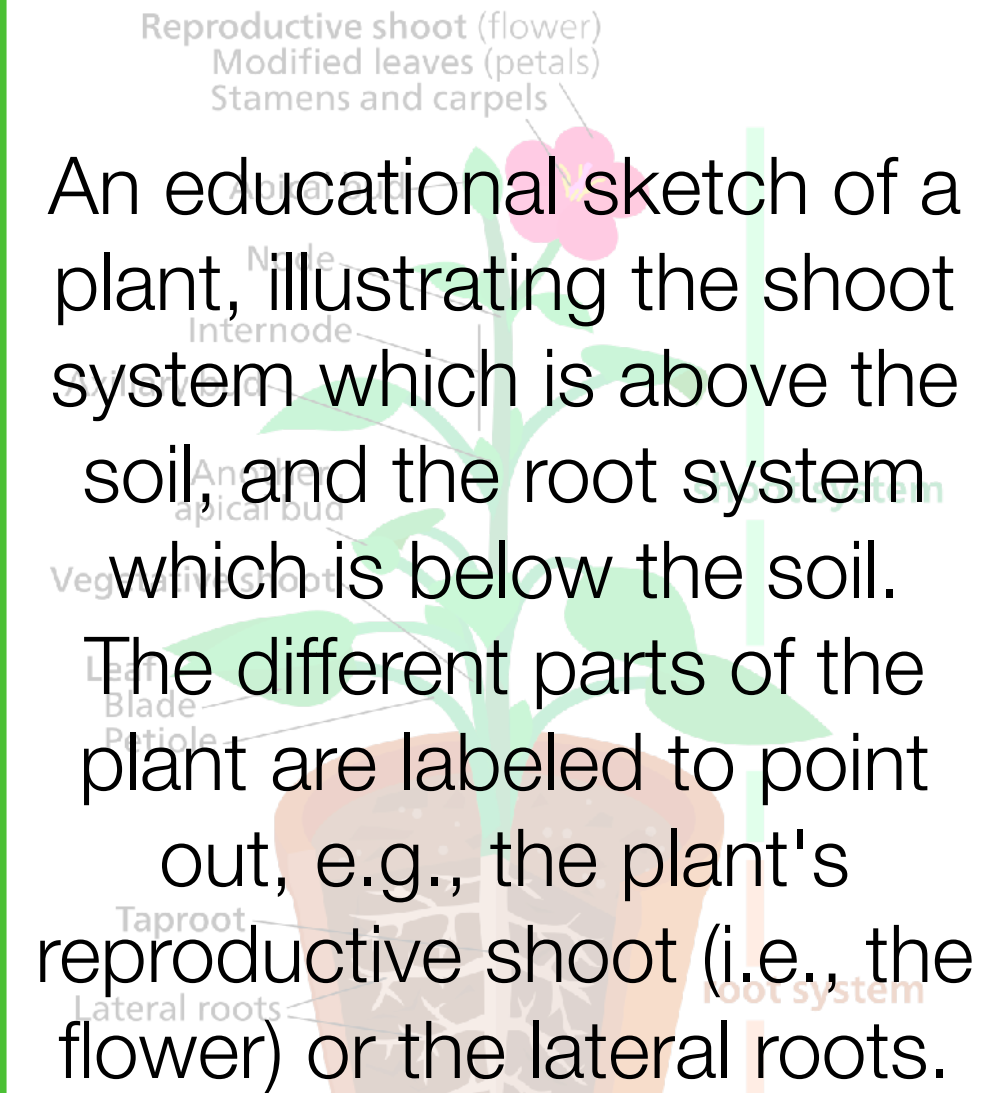
An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.

Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.

Description
accessibility
replaces image

Caption
complements image

The Communicative Goal of the Image-Based **Text**



An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.

Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.

Description
accessibility
replaces image

Captions and descriptions are both image-based, but are produced in complementary **settings**.

(Goodwin & Duranti; 1992)

Caption
complements image

The Communicative Goal of the Image-Based **Text**

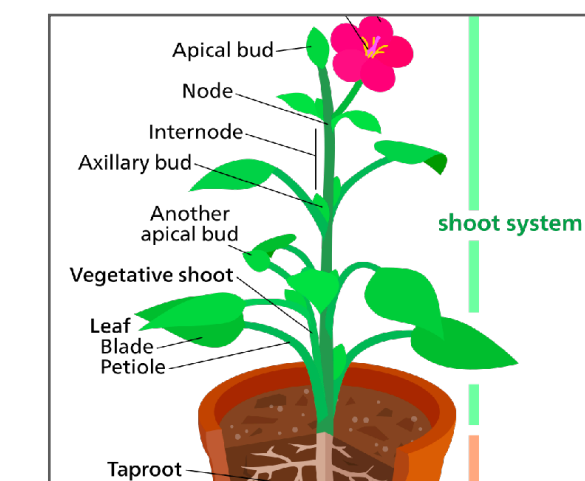
In AI, we tend to reduce them to the same problem.

Description
accessibility
replaces image

Caption
complements image

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.

Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.



Model

Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.

Image-based text generation depends on ...

1 the **image-based text**'s communicative goal.

→ description \neq caption

2 the **image**'s communicative goal.

A sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. Part of the root system are the taproots and lateral roots. The taproot refers to the central root and the lateral roots are the smaller side roots that ...

A diagram of the anatomy of a plant with labels of structural parts of the plant and the roots.

Finding the Image-Based **Text**'s Communicative Goal

Concadia: A naturalistic image-based text dataset from Wikipedia

96,918 images with captions, alt descriptions and surrounding paragraph

Finding the Image-Based **Text**'s Communicative Goal

Concadia: A naturalistic image-based text dataset from Wikipedia

96,918 images with captions, alt descriptions and surrounding paragraph

Wikipedia-Article on **Banana**

image context: In global commerce in 2009, by far the most important cultivars belonged to the triploid AAA group of *Musa acuminata*, commonly referred to as Cavendish group bananas. They accounted for the majority of banana exports, despite only coming into existence in 1836. The cultivars Dwarf Cavendish and Grand Nain (Chiquita Banana) gained popularity in the 1950s after the previous mass-produced cultivar, Gros Michel (also an AAA ...



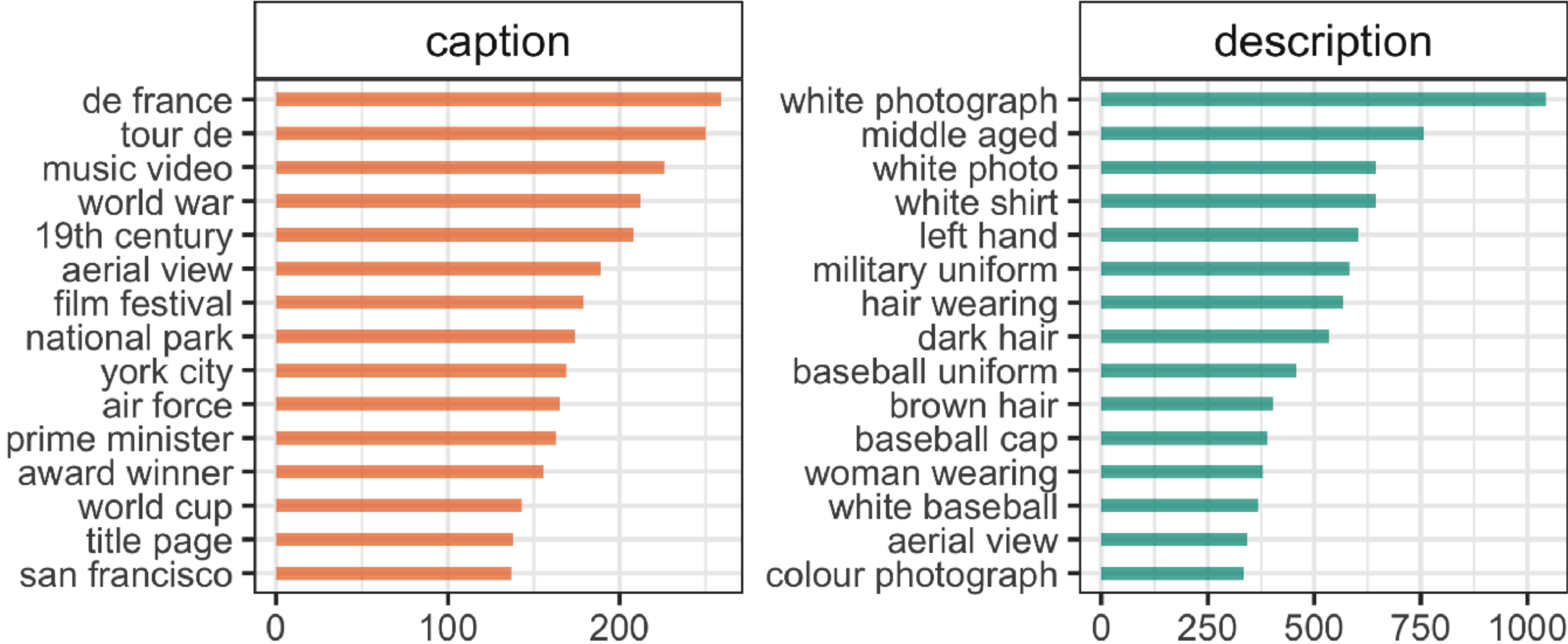
(Accessibility) **Description:**

Grocery store photo of several bunches of bananas

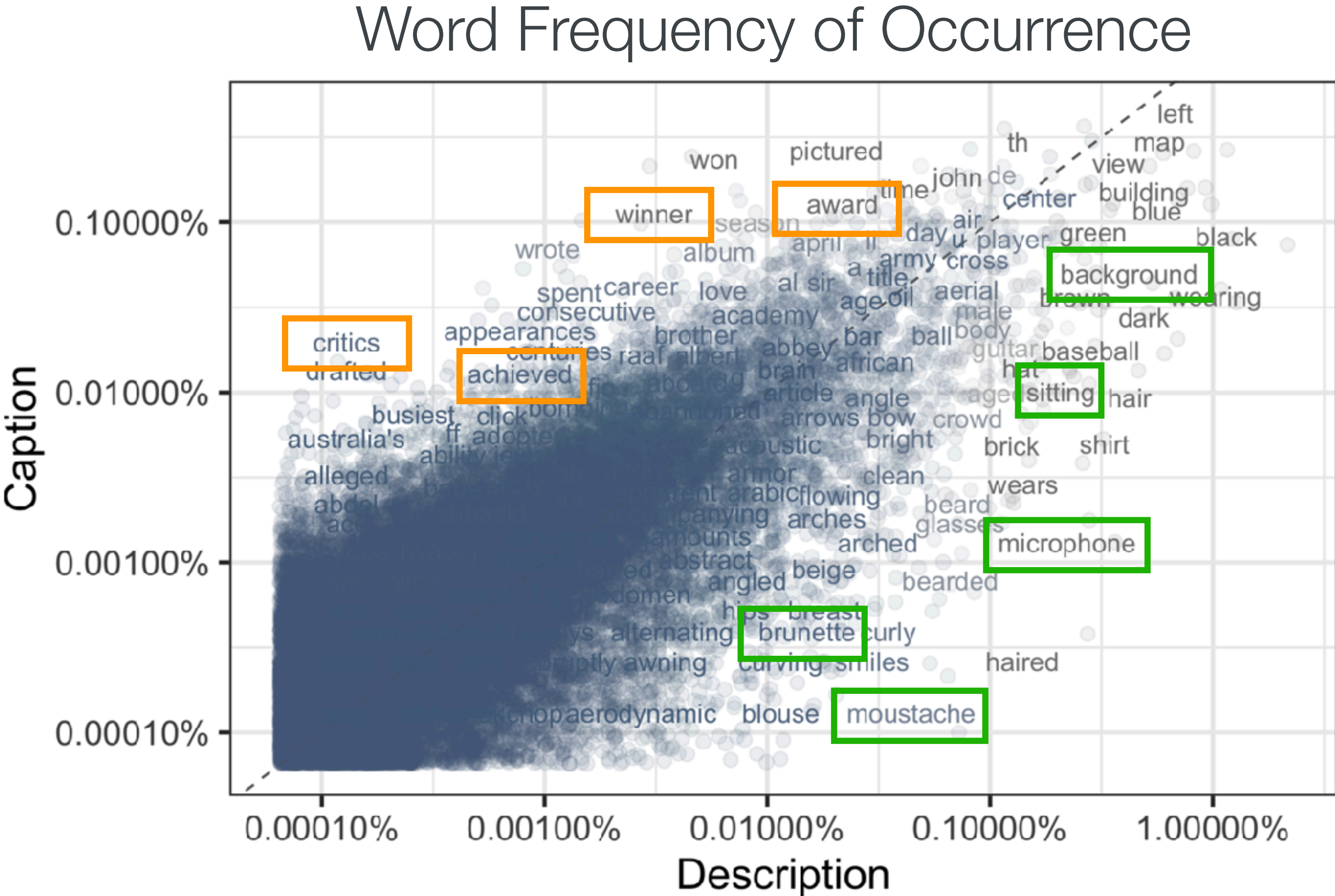
(Contextualizing) **Caption:**

Cavendish bananas are the main commercial banana cultivars sold in the world market.

Concadia: A naturalistic image-based text dataset

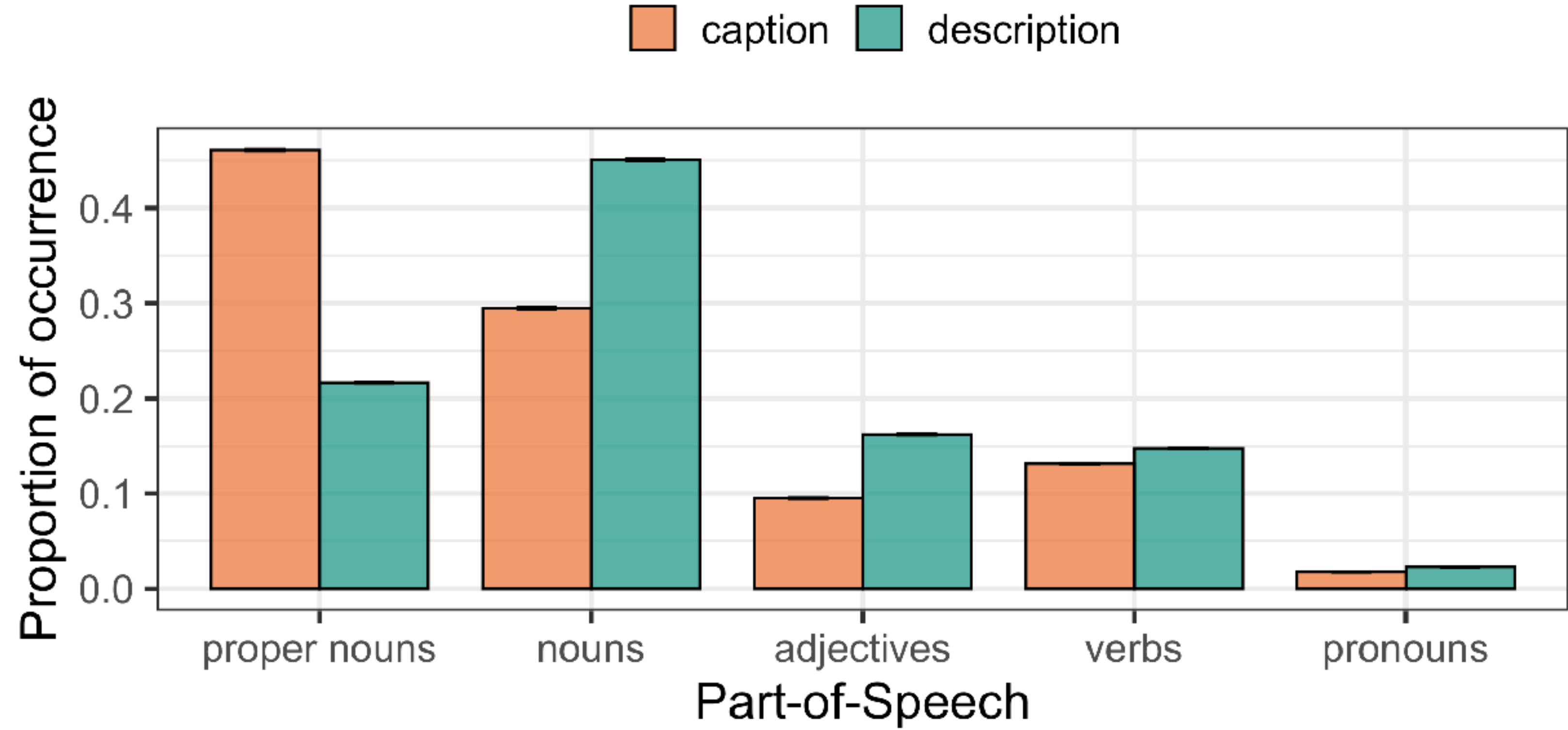


Concadia: A naturalistic image-based text dataset



Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Concadia: A naturalistic image-based text dataset

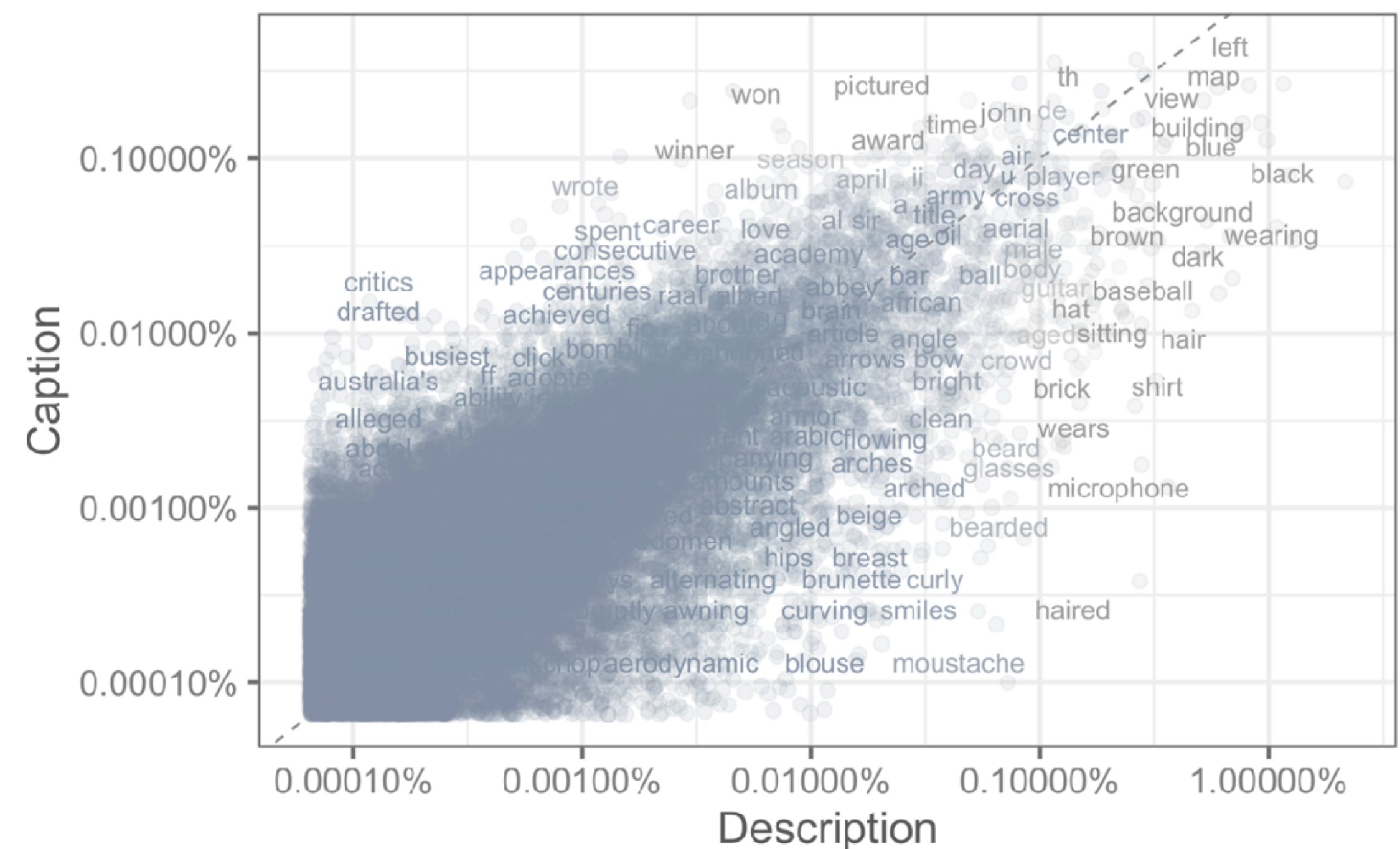
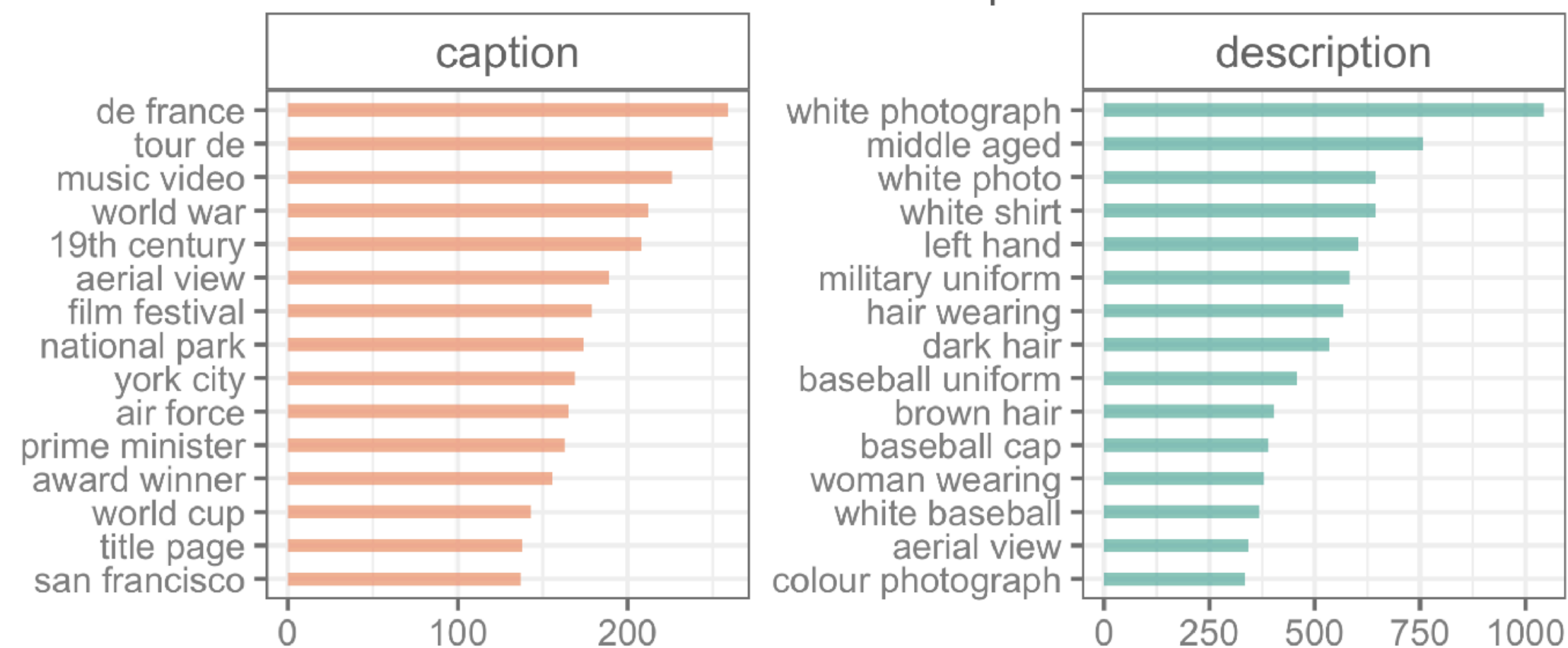
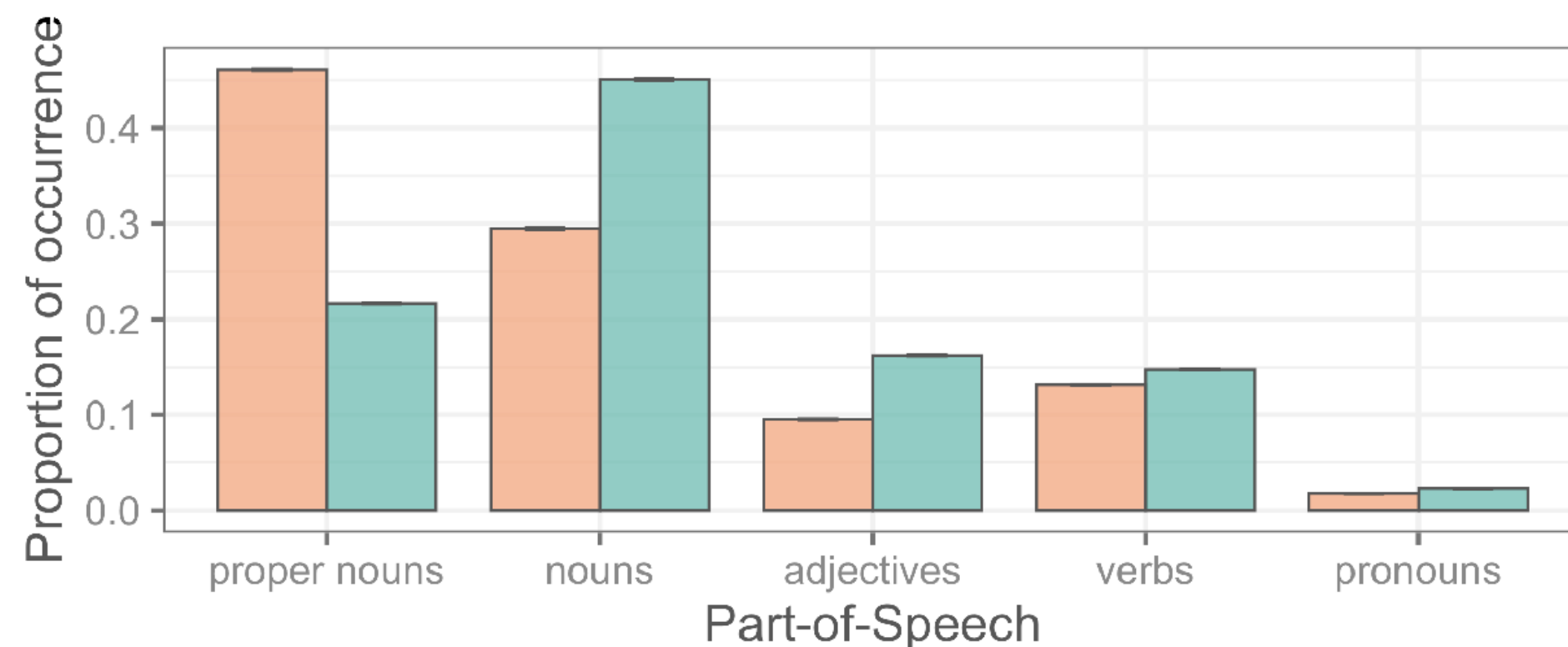


The Image-Based **Text**'s Communicative Goal



Large-scale analysis of naturalistic data:

The content of alt descriptions and captions differs in structured ways.



Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

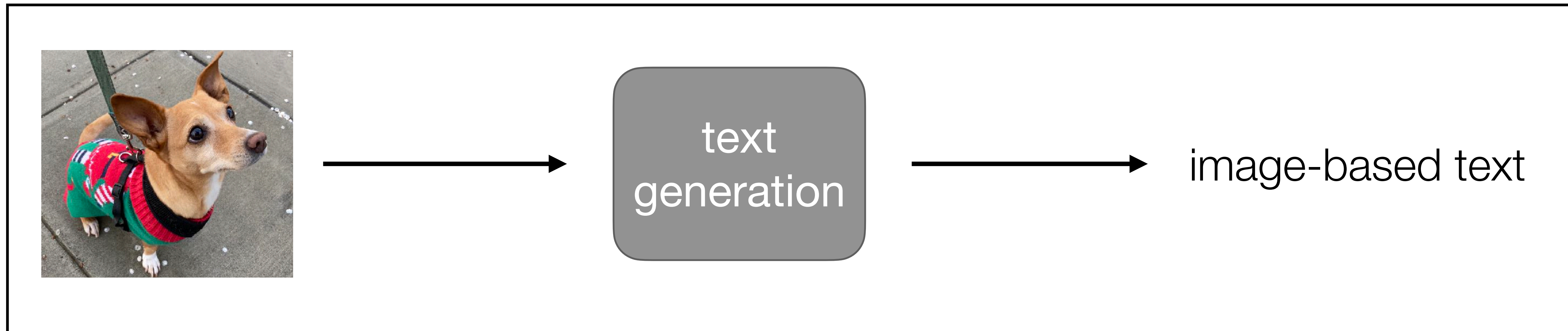
Generating Image-Based Text with a Purpose

Architecture

ResNet + LSTM

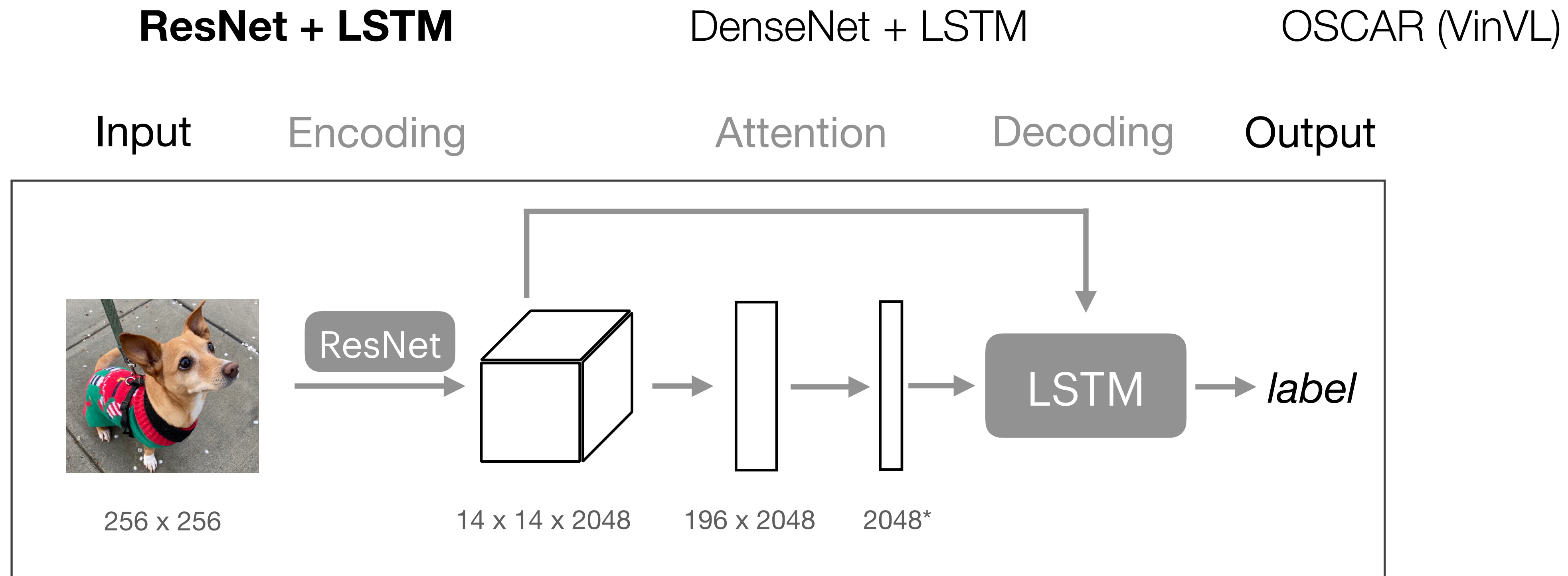
DenseNet + LSTM

OSCAR (VinVL)



Generating Image-Based Text with a Purpose

Architecture

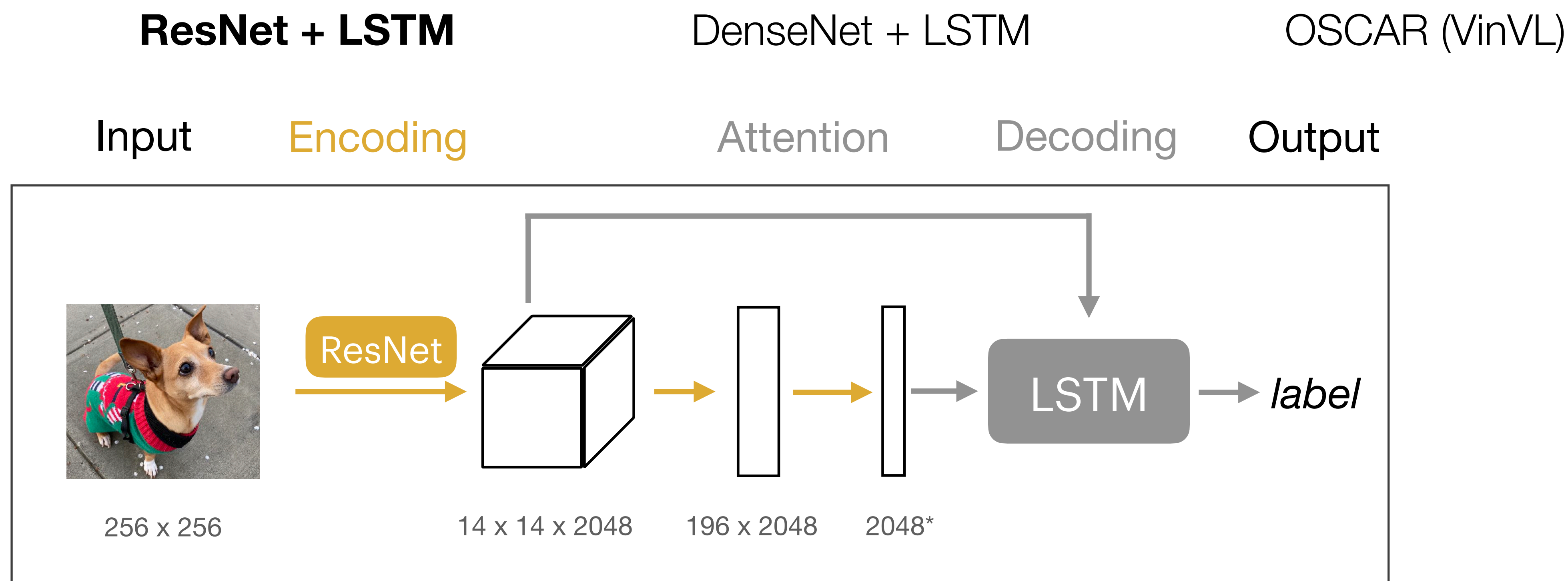


*simplified but will be expanded later in the talk

Based on Xu et al., 2015, He et al., 2016

Generating Image-Based Text with a Purpose

Architecture

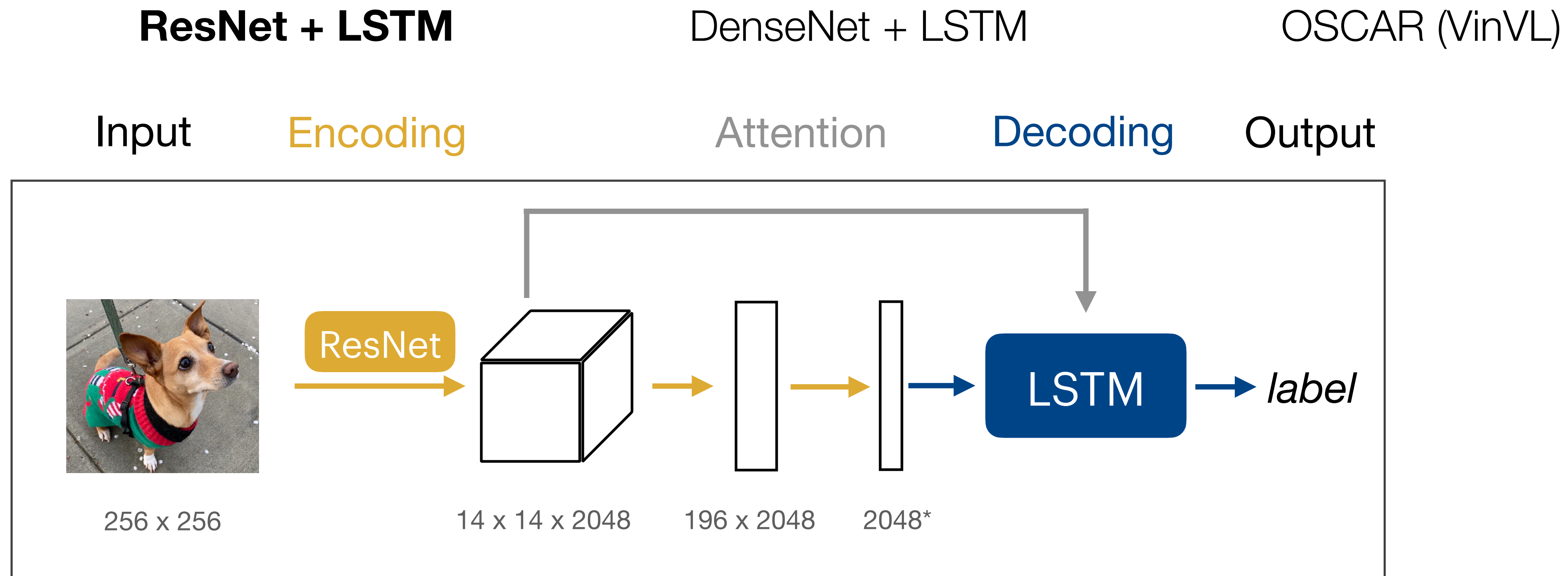


*simplified but will be expanded later in the talk

Based on Xu et al., 2015, He et al., 2016

Generating Image-Based Text with a Purpose

Architecture

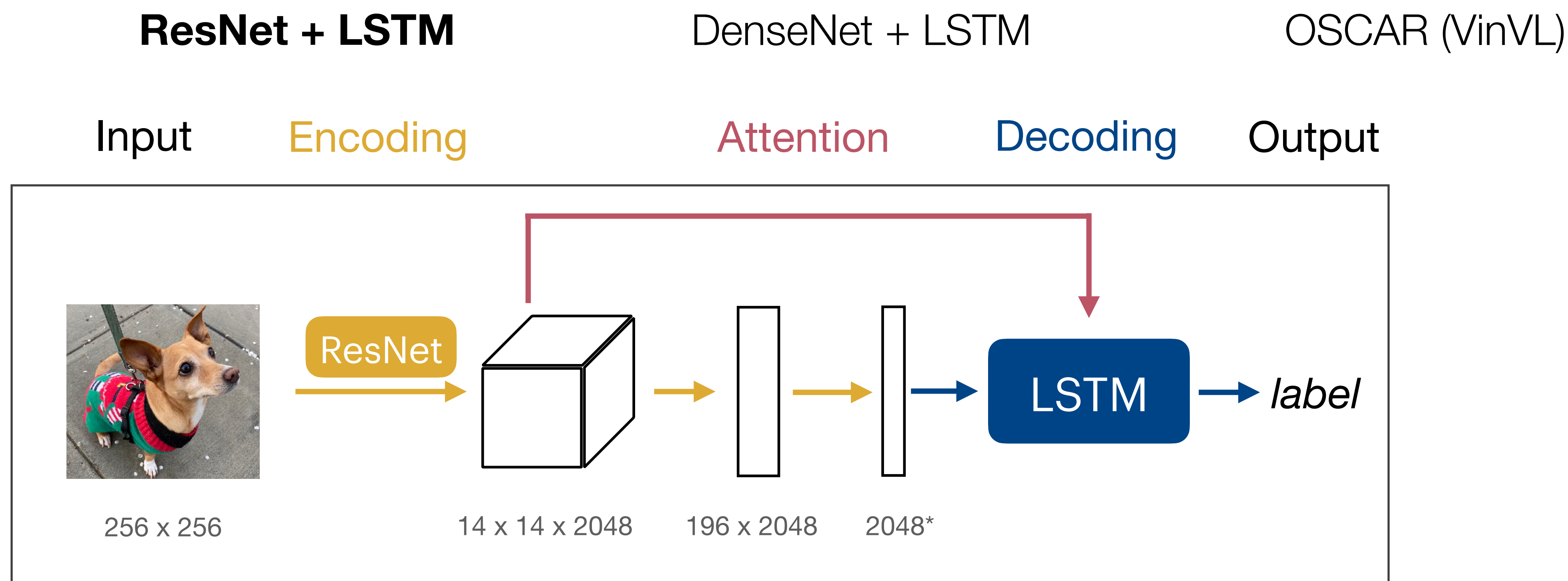


*simplified but will be expanded later in the talk

Based on Xu et al., 2015, He et al., 2016

Generating Image-Based Text with a Purpose

Architecture

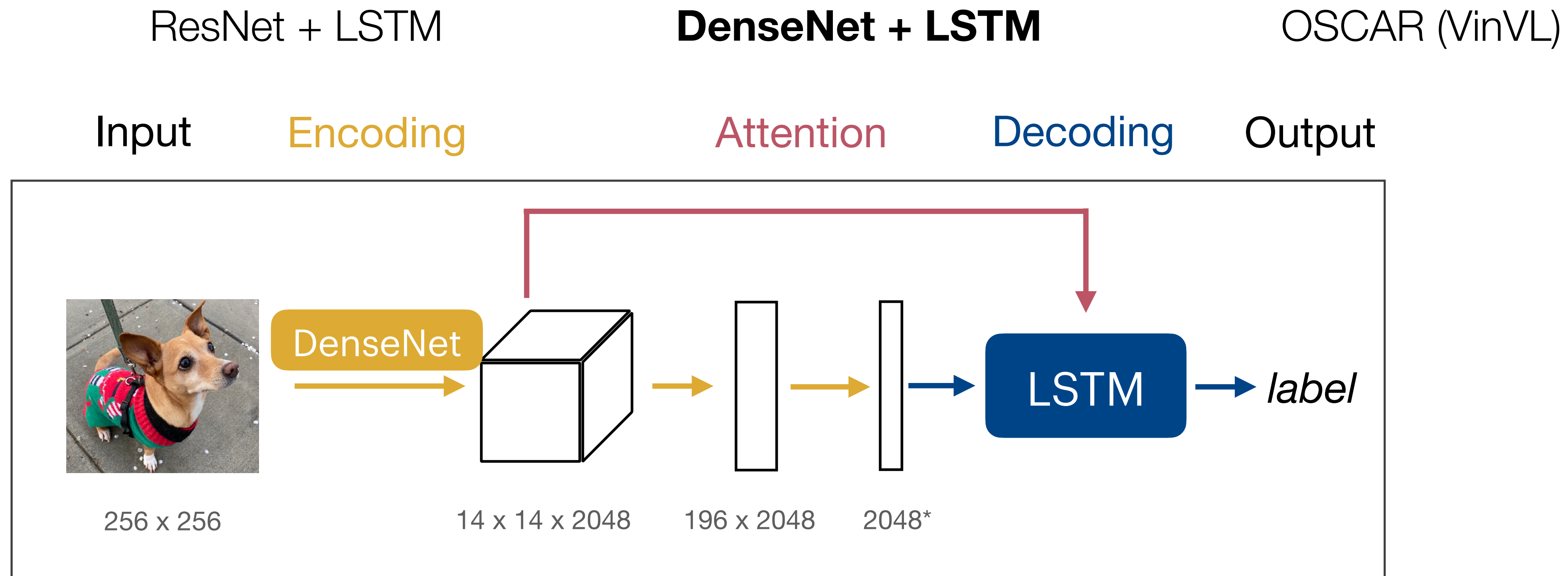


*simplified but will be expanded later in the talk

Based on Xu et al., 2015, He et al., 2016

Generating Image-Based Text with a Purpose

Architecture

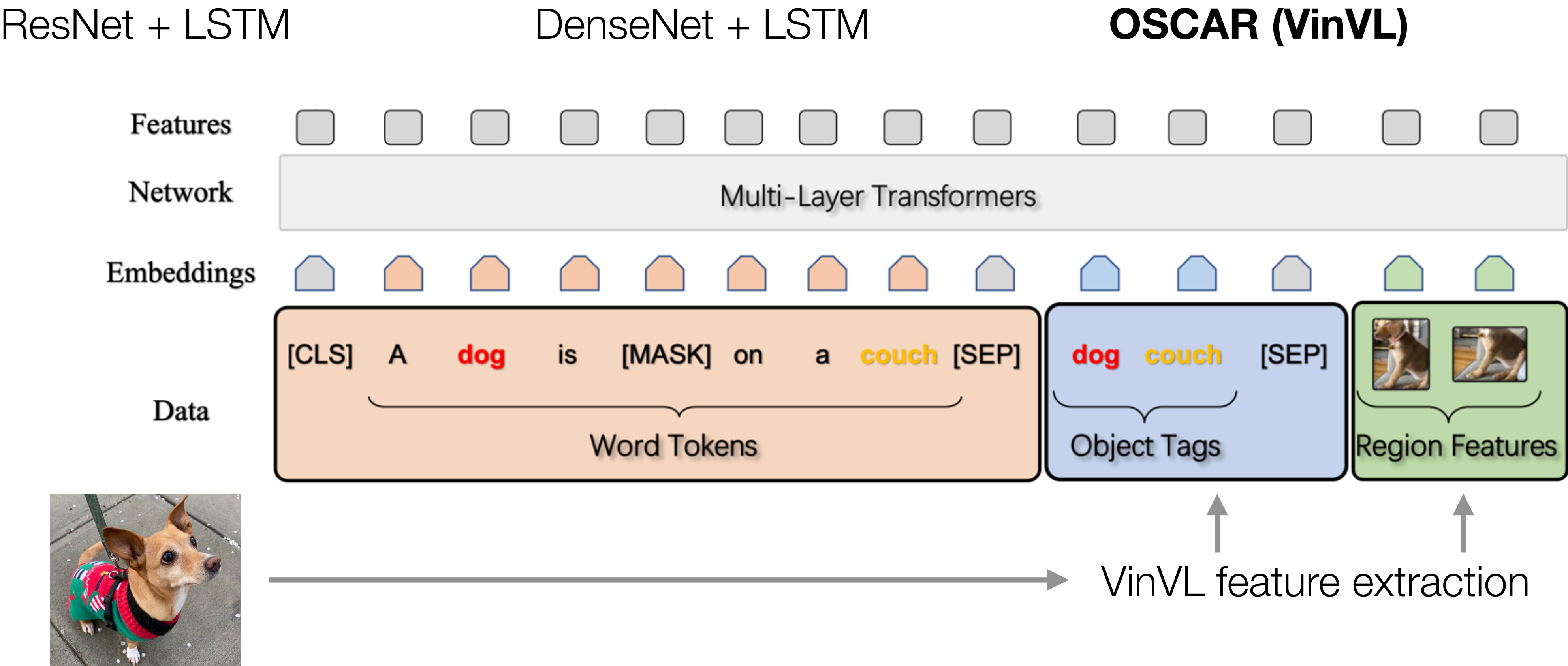


*simplified but will be expanded later in the talk

Based on Deng et al., 2020; Hossain et al., 2021

Generating Image-Based Text with a Purpose

Architecture



Generating Image-Based Text with a Purpose

Training

ResNet + LSTM

DenseNet + LSTM

OSCAR (VinVL)

1

2

Generating Image-Based Text with a Purpose

Training

ResNet + LSTM

DenseNet + LSTM

OSCAR (VinVL)

Caption generation

1



image pre-
processing



text
generation



Cavendish bananas are
the main commercial
banana cultivar sold in the
world market.

2

Generating Image-Based Text with a Purpose

Training

ResNet + LSTM

DenseNet + LSTM

OSCAR (VinVL)

1

Caption generation



image pre-
processing



text
generation



Cavendish bananas are
the main commercial
banana cultivar sold in the
world market.

2

Description generation



image pre-
processing

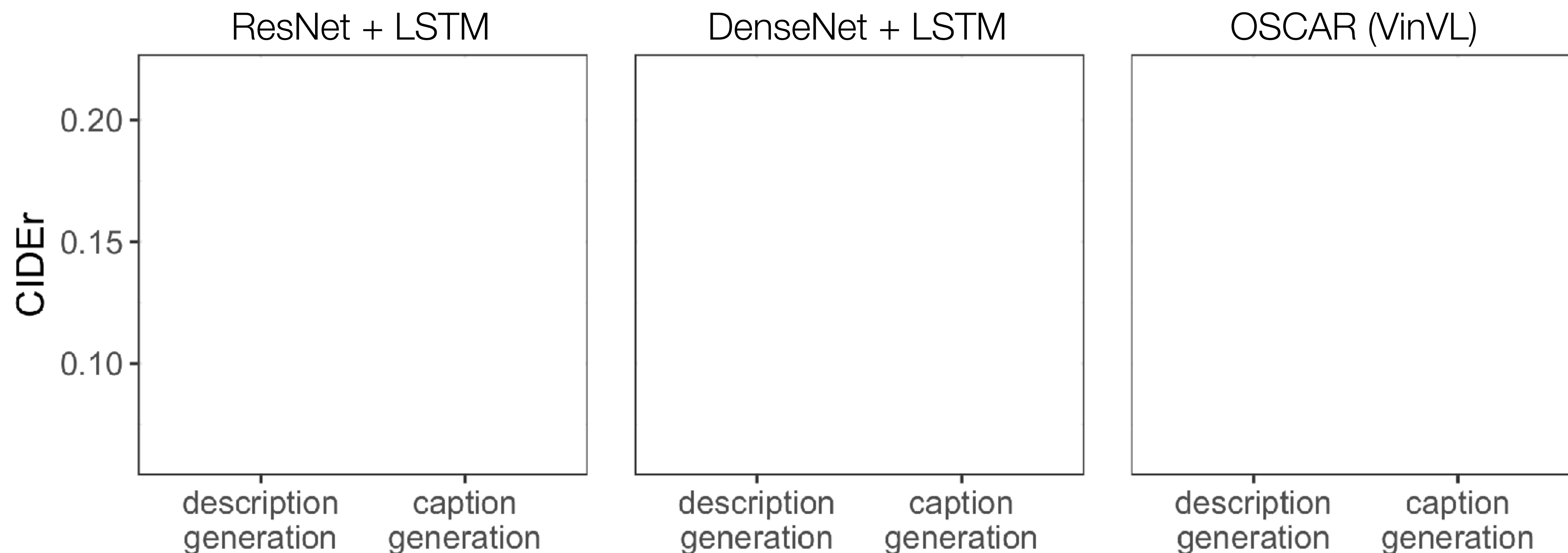


text
generation



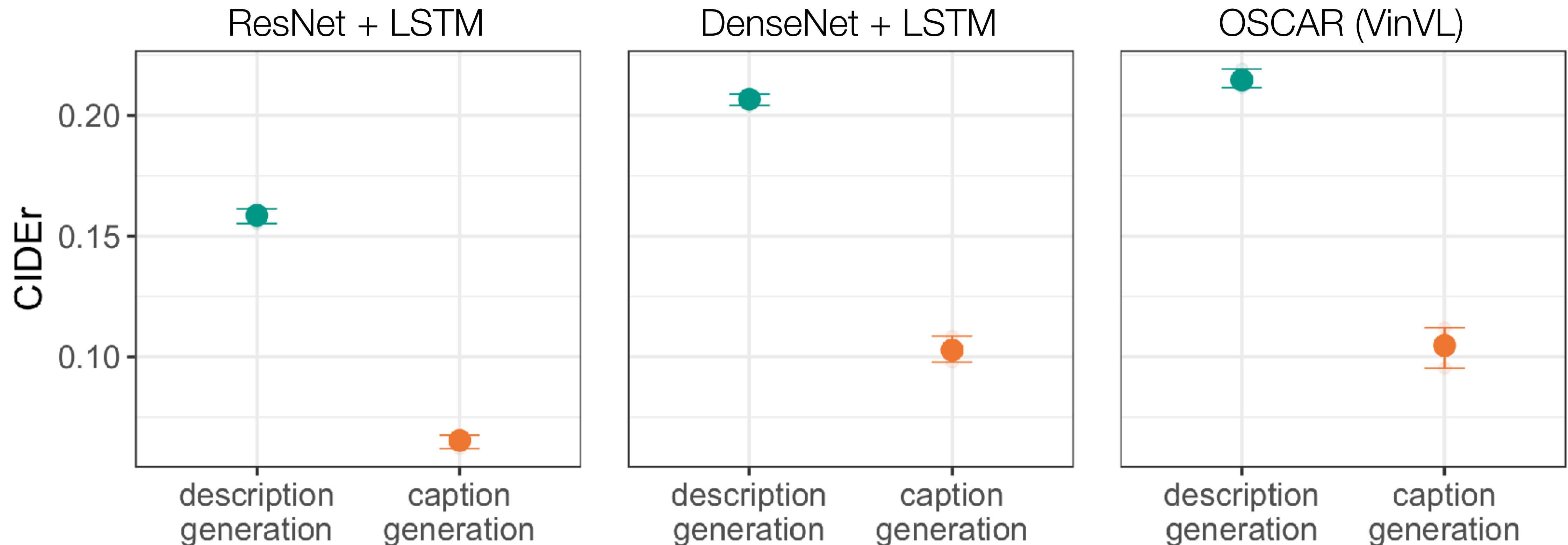
Grocery store photo of
several bunches of
banana

Generating Image-Based Text with a Purpose



Generating Image-Based Text with a Purpose

Caption generation is harder than description generation



The Image-Based **Text**'s Communicative Goal



Large-scale analysis of naturalistic data:

The content of alt descriptions and captions differs in structured ways.



Models of image-based text generation:

Learning alt description and caption generation are distinct challenges.

The Image-Based **Text**'s Communicative Goal



Large-scale analysis of naturalistic data:

The content of alt descriptions and captions differs in structured ways.



Models of image-based text generation:

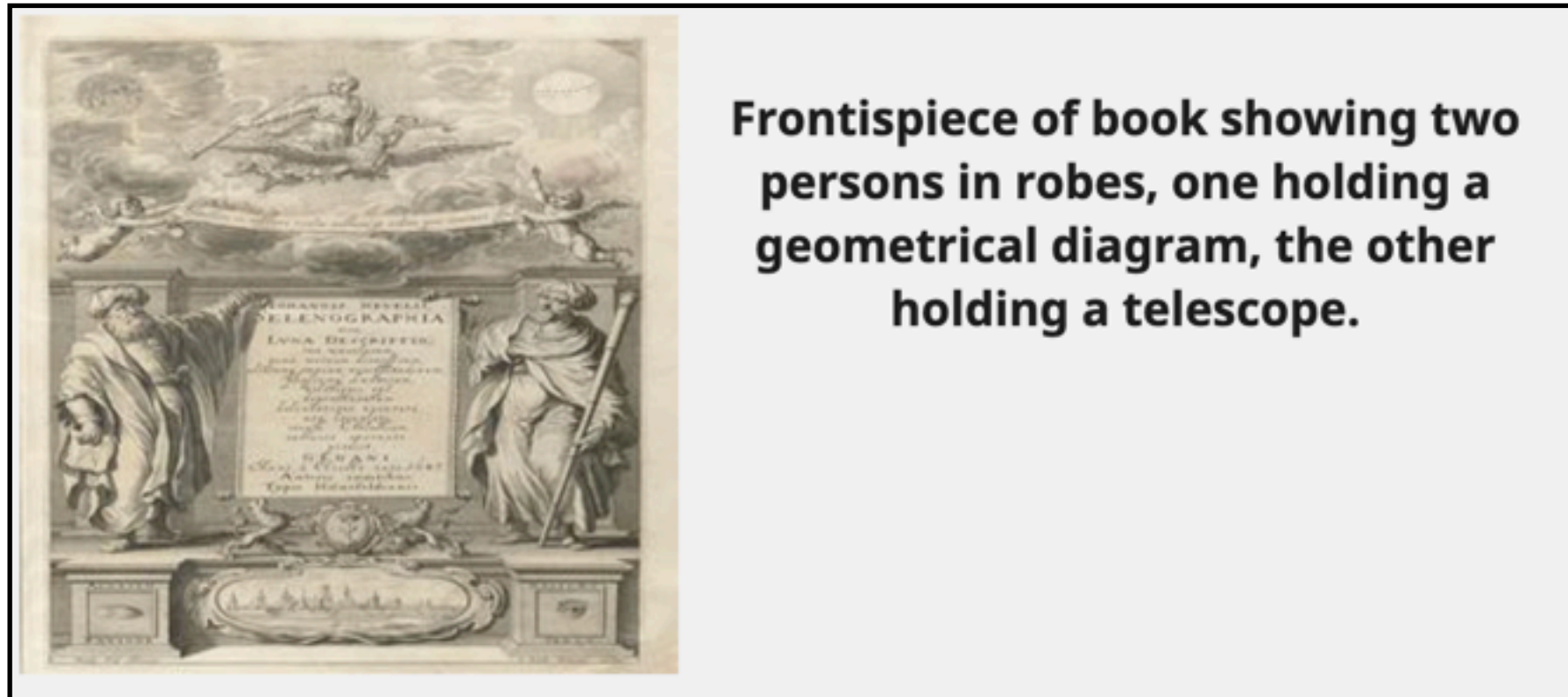
Learning alt description and caption generation are distinct challenges.

Is the distinction between descriptions and captions due to a distinction in the text's **communicative goal**?

Is this distinction reflected in **models** trained on the respective data?

Evaluating the Purpose of Descriptions & Captions

preregistered human subject experiment




Frontispiece of book showing two persons in robes, one holding a geometrical diagram, the other holding a telescope.

→ image,
text (description: human/model,
caption: human/model)

Evaluating the Purpose of Descriptions & Captions

preregistered human subject experiment



Frontispiece of book showing two persons in robes, one holding a geometrical diagram, the other holding a telescope.

Q1: How **useful** would the **text alone** be to help someone imagine this picture (e.g, a visually impaired person)?

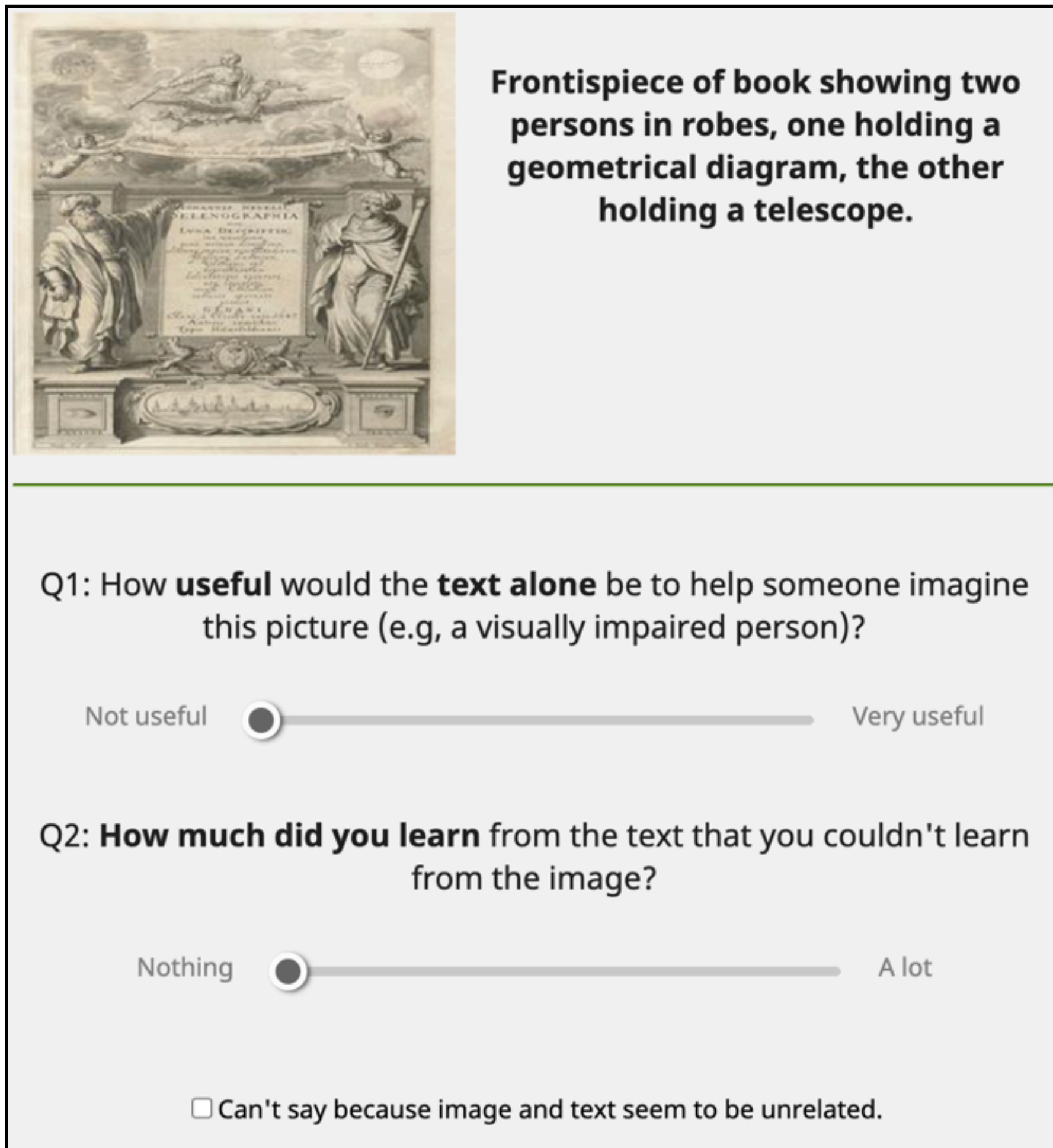
Not useful Very useful

→ image,
text (description: human/model,
caption: human/model)

→ question for descriptive quality

Evaluating the Purpose of Descriptions & Captions

preregistered human subject experiment




→ image,
text (description: human/model,
caption: human/model)

→ question for descriptive quality

→ question for caption quality

Towards AI for Image Accessibility

preregistered human subject experiment



Frontispiece of book showing two persons in robes, one holding a geometrical diagram, the other holding a telescope.

Q1: How **useful** would the **text alone** be to help someone imagine this picture (e.g, a visually impaired person)?

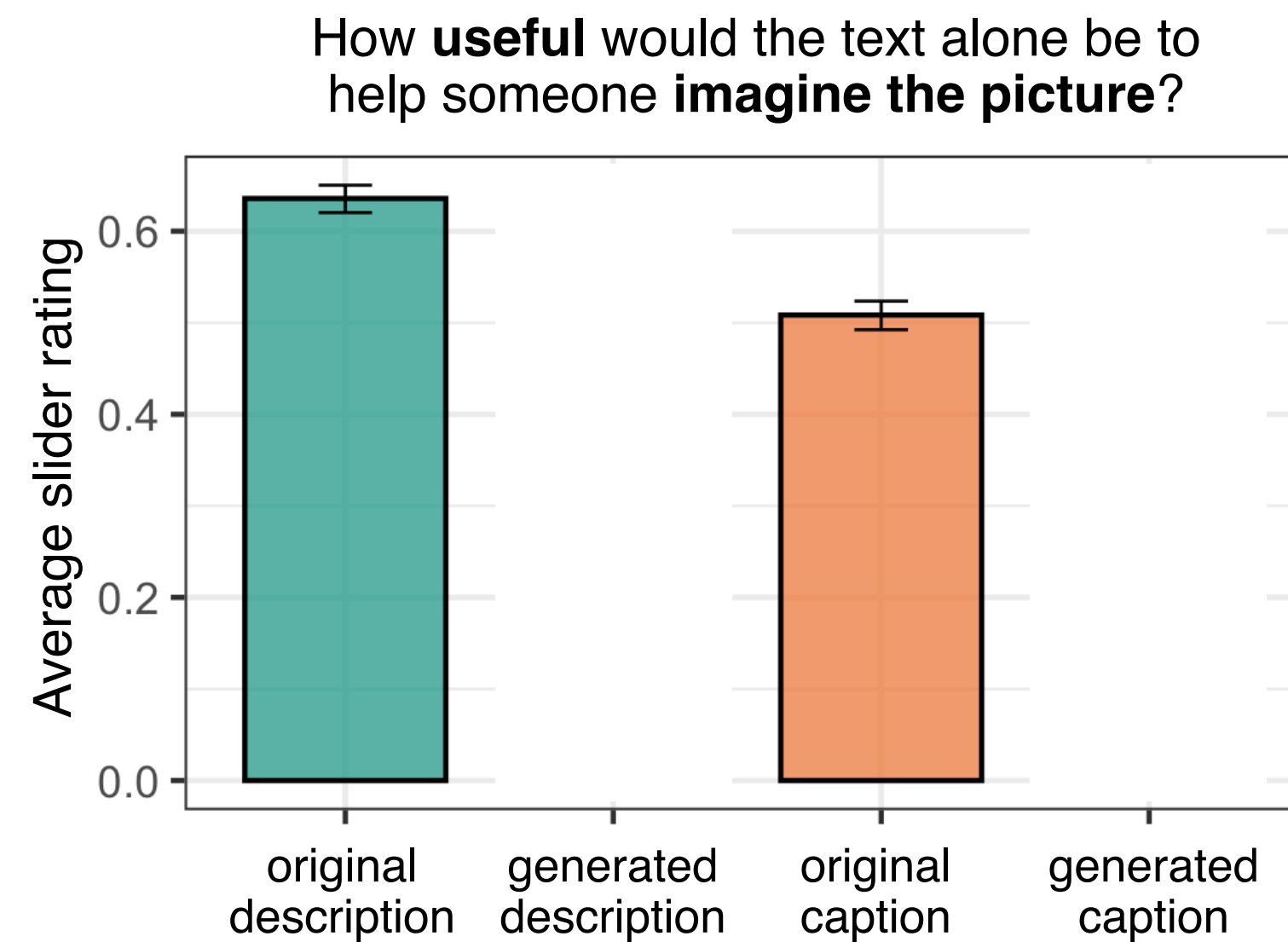
Not useful ☐ ☐ ☐ ☐ ☐ Very useful

Q2: How much did you learn from the text that you couldn't learn from the image?

Nothing ☐ ☐ ☐ ☐ A lot

☐ Can't say because image and text seem to be unrelated.


question for descriptive quality



Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Towards AI for Image Accessibility

preregistered human subject experiment



Frontispiece of book showing two persons in robes, one holding a geometrical diagram, the other holding a telescope.

Q1: How **useful** would the **text alone** be to help someone imagine this picture (e.g, a visually impaired person)?

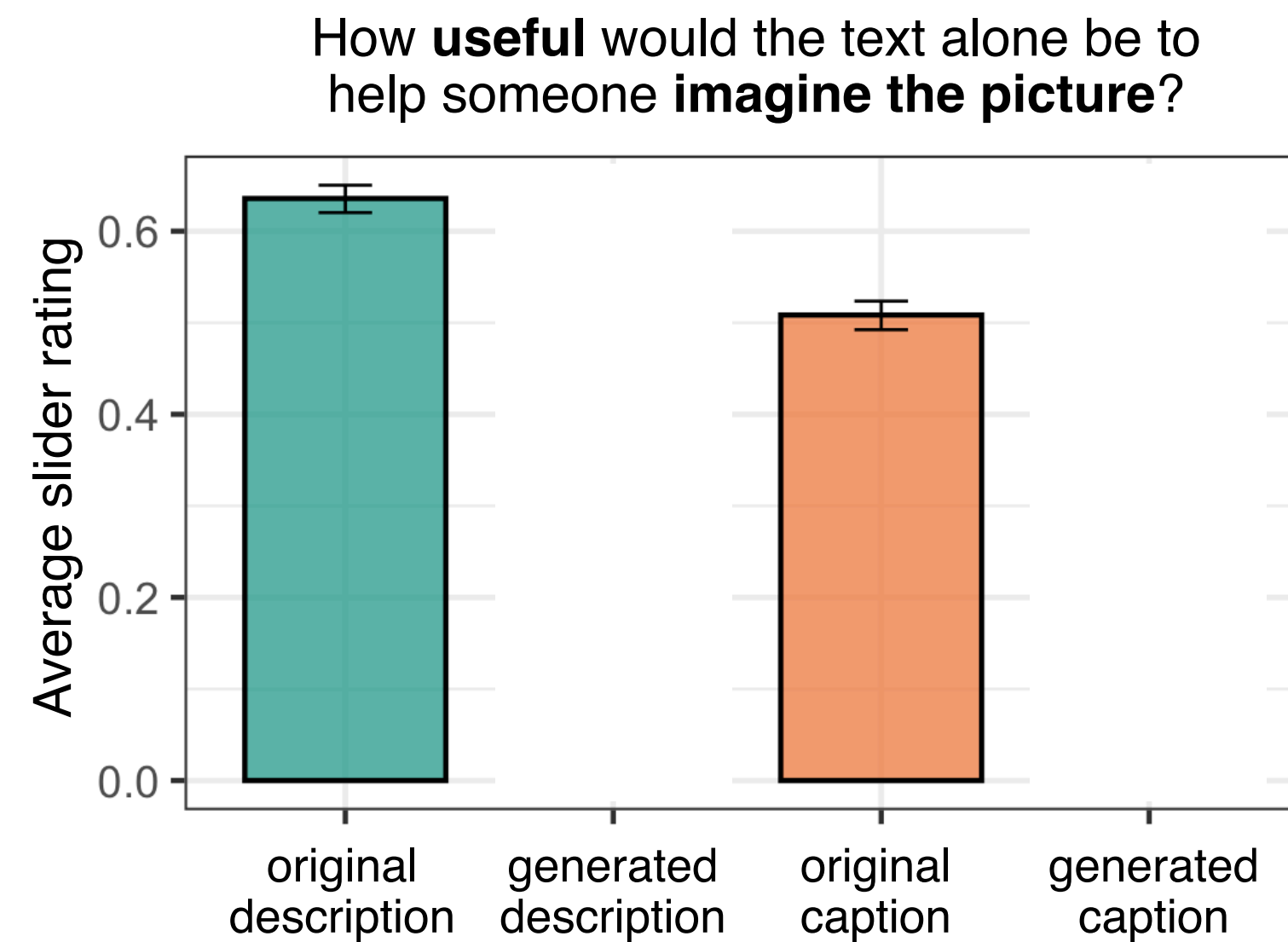
Not useful Very useful

Q2: How much did you **learn** from the text that you couldn't learn from the image?

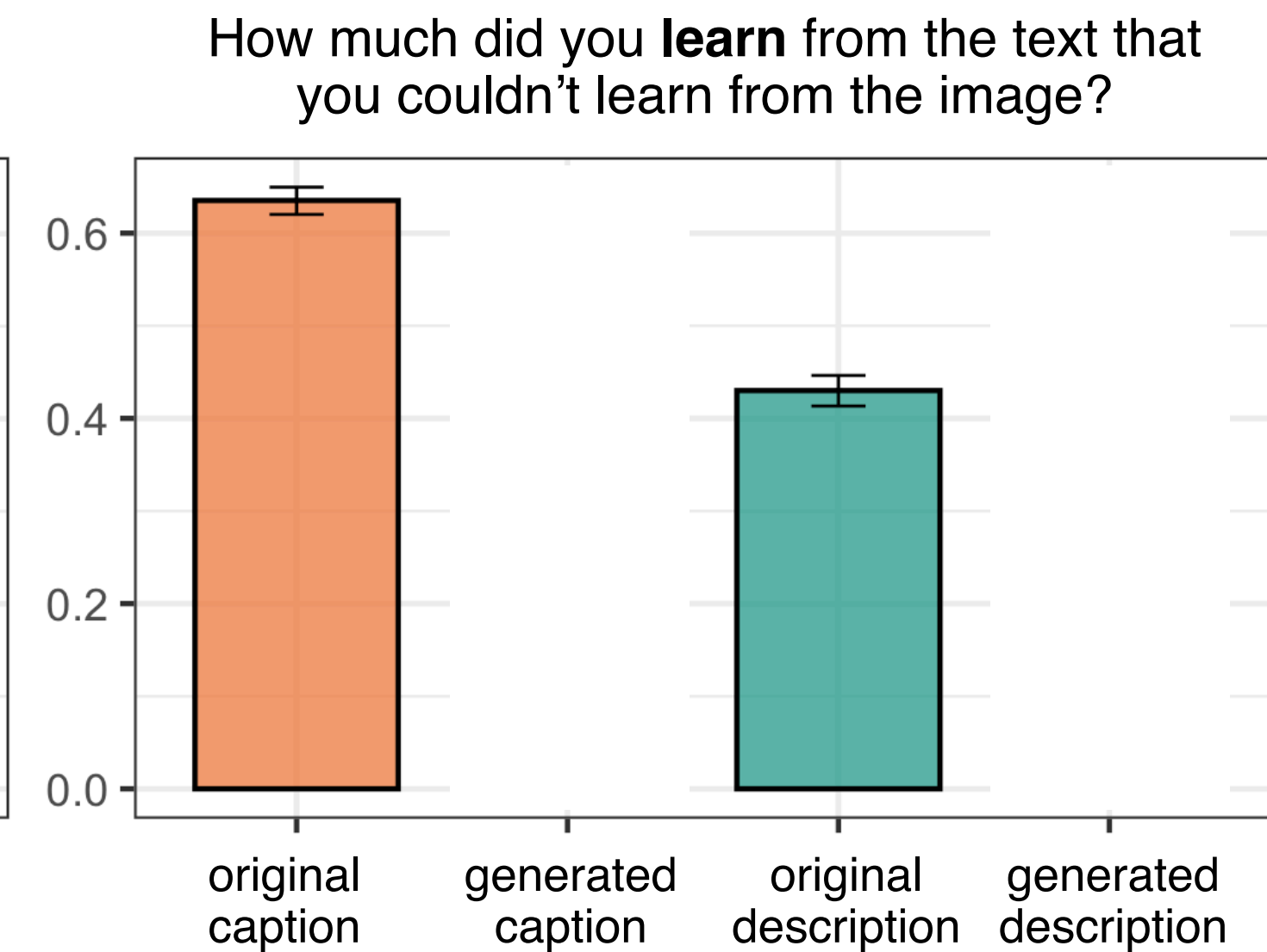
Nothing A lot

☐ Can't say because image and text seem to be unrelated.

question for descriptive quality




question for caption quality



Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Towards AI for Image Accessibility

preregistered human subject experiment



Frontispiece of book showing two persons in robes, one holding a geometrical diagram, the other holding a telescope.

Q1: How **useful** would the **text alone** be to help someone imagine this picture (e.g, a visually impaired person)?

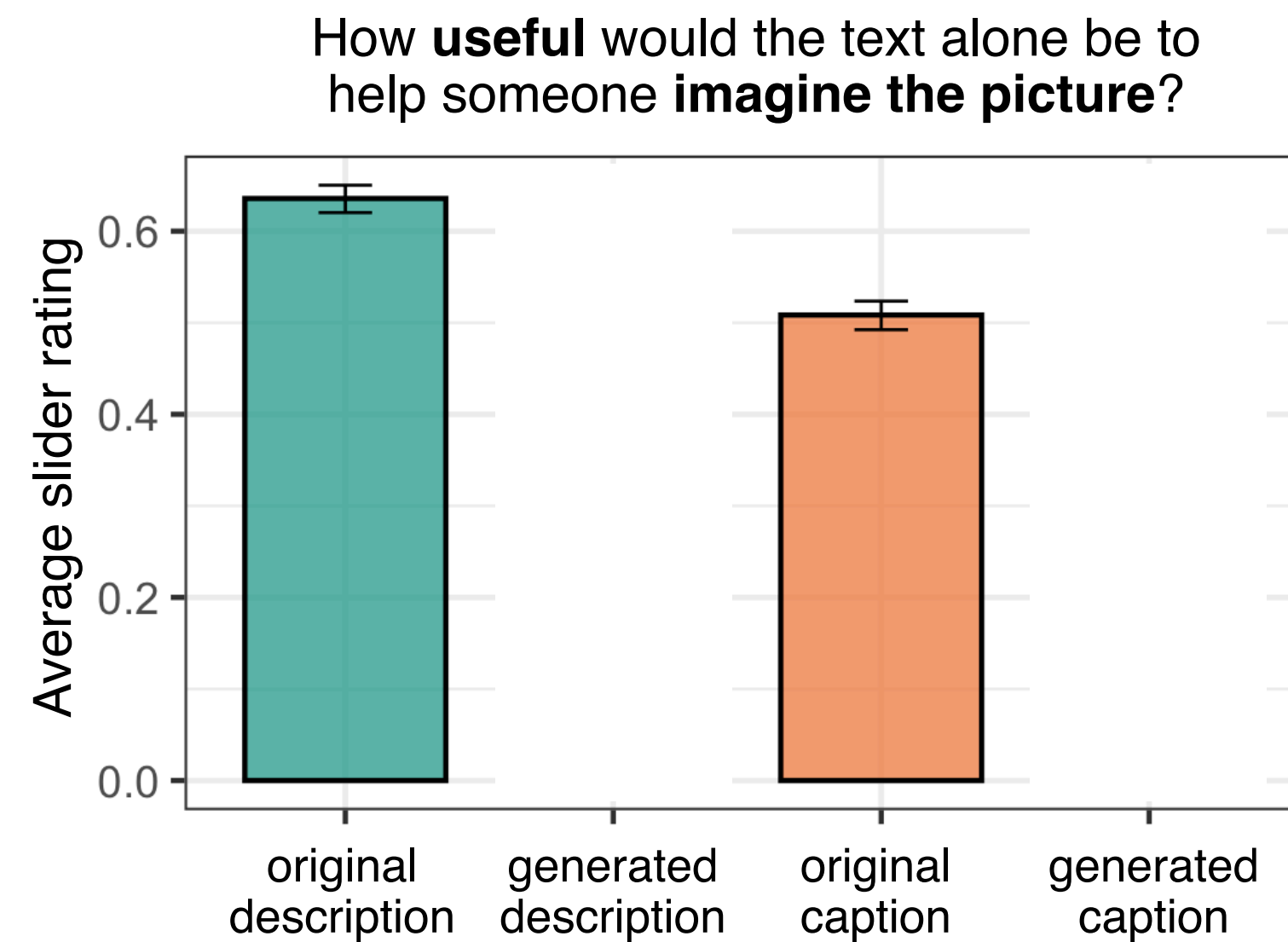
Not useful Very useful

Q2: How much did you **learn** from the text that you couldn't learn from the image?

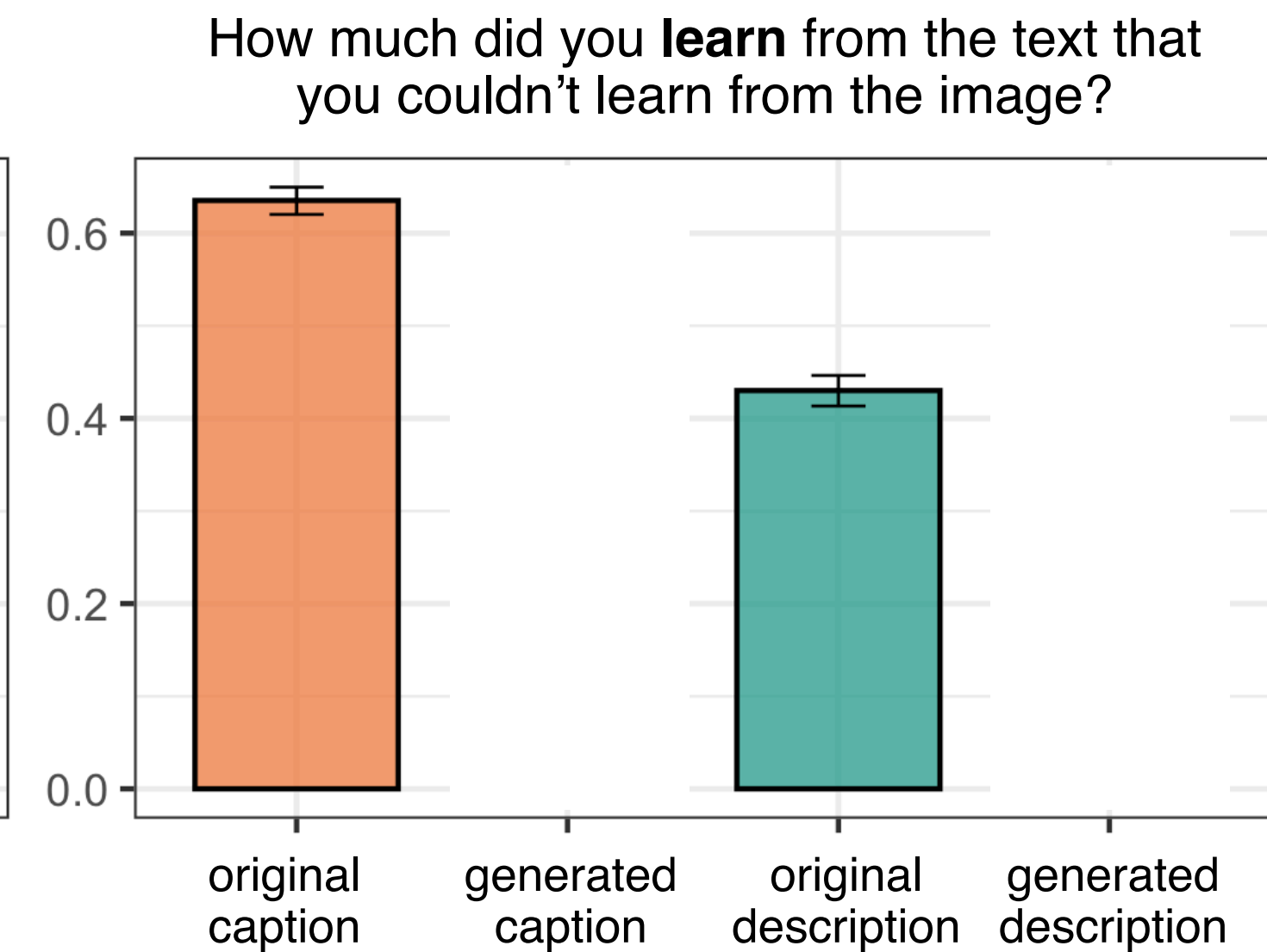
Nothing A lot

☐ Can't say because image and text seem to be unrelated.

question for descriptive quality



question for caption quality




Descriptions and captions optimize for distinct communicative goals.

Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Towards AI for Image Accessibility

preregistered human subject experiment



Frontispiece of book showing two persons in robes, one holding a geometrical diagram, the other holding a telescope.

Q1: How **useful** would the **text alone** be to help someone imagine this picture (e.g, a visually impaired person)?

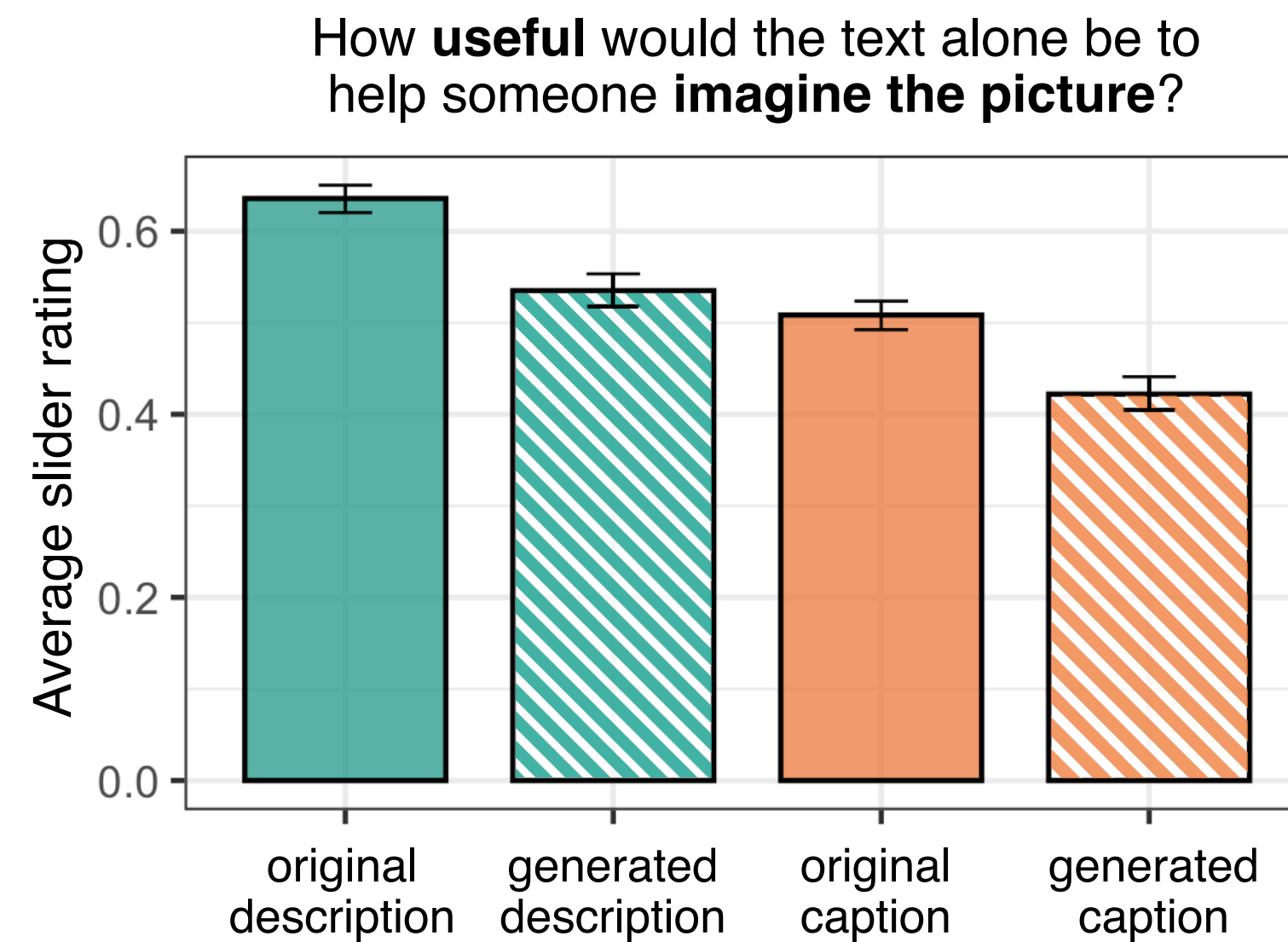
Not useful Very useful

Q2: How much did you **learn** from the text that you couldn't learn from the image?

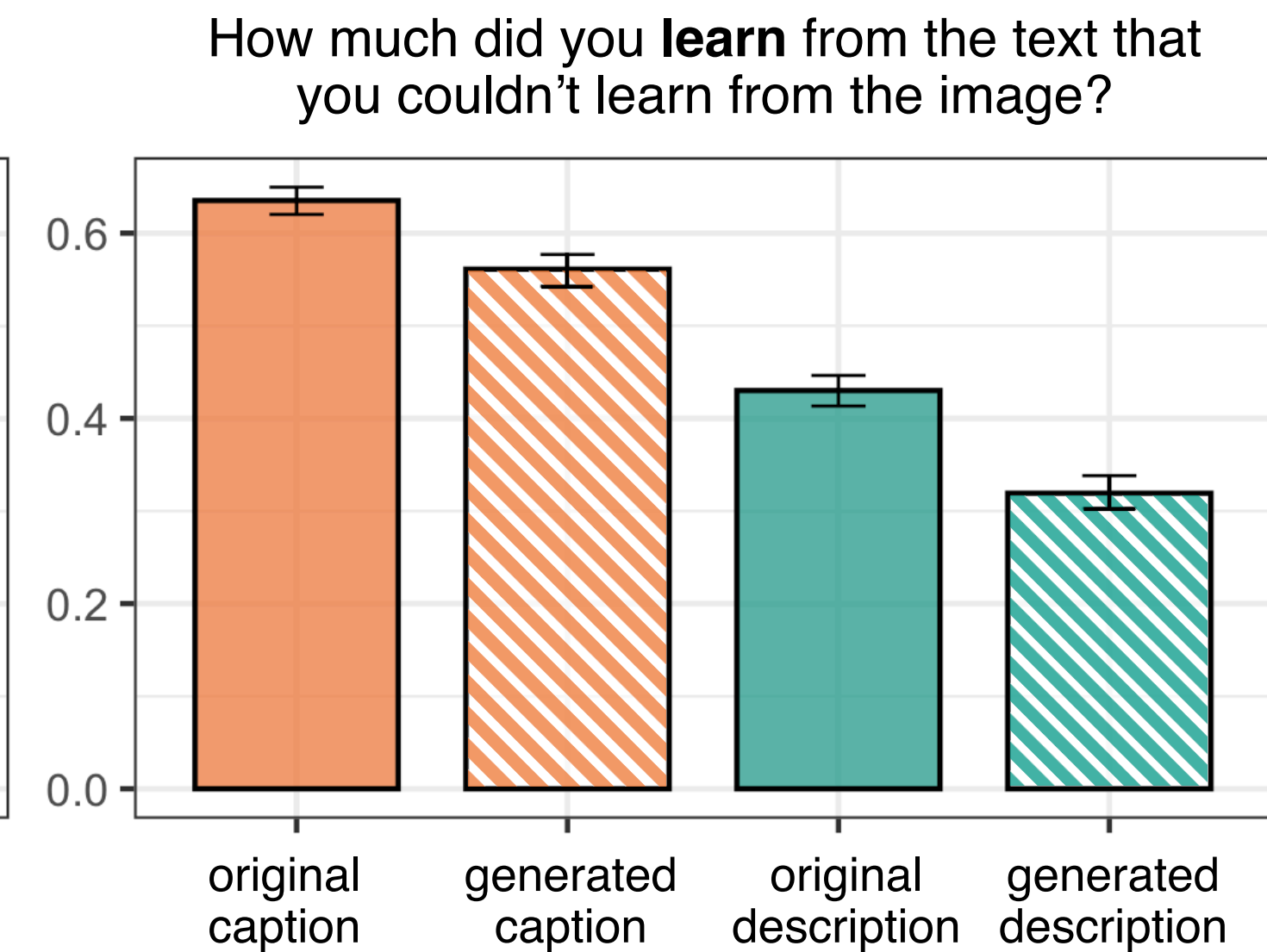
Nothing A lot

☐ Can't say because image and text seem to be unrelated.

question for descriptive quality



question for caption quality



Descriptions and captions optimize for distinct communicative goals.

Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

The Image-Based **Text**'s Communicative Goal



Large-scale analysis of naturalistic data:

The content of alt descriptions and captions differs in structured ways.



Models of image-based text generation:

Learning alt description and caption generation are distinct challenges.



Human experiment:

Descriptions and captions optimize for distinct communicative goals.

This distinction is reflected in models trained on the respective data.

Image-based text generation depends on ...

1 the **image-based text**'s communicative goal.

→ description \neq caption

2 the **image**'s communicative goal.

A sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. Part of the root system are the taproots and lateral roots. The taproot refers to the central root and the lateral roots are the smaller side roots that ...

A diagram of the anatomy of a plant with labels of structural parts of the plant and the roots.

A Pragmatic Approach to Image Descriptions

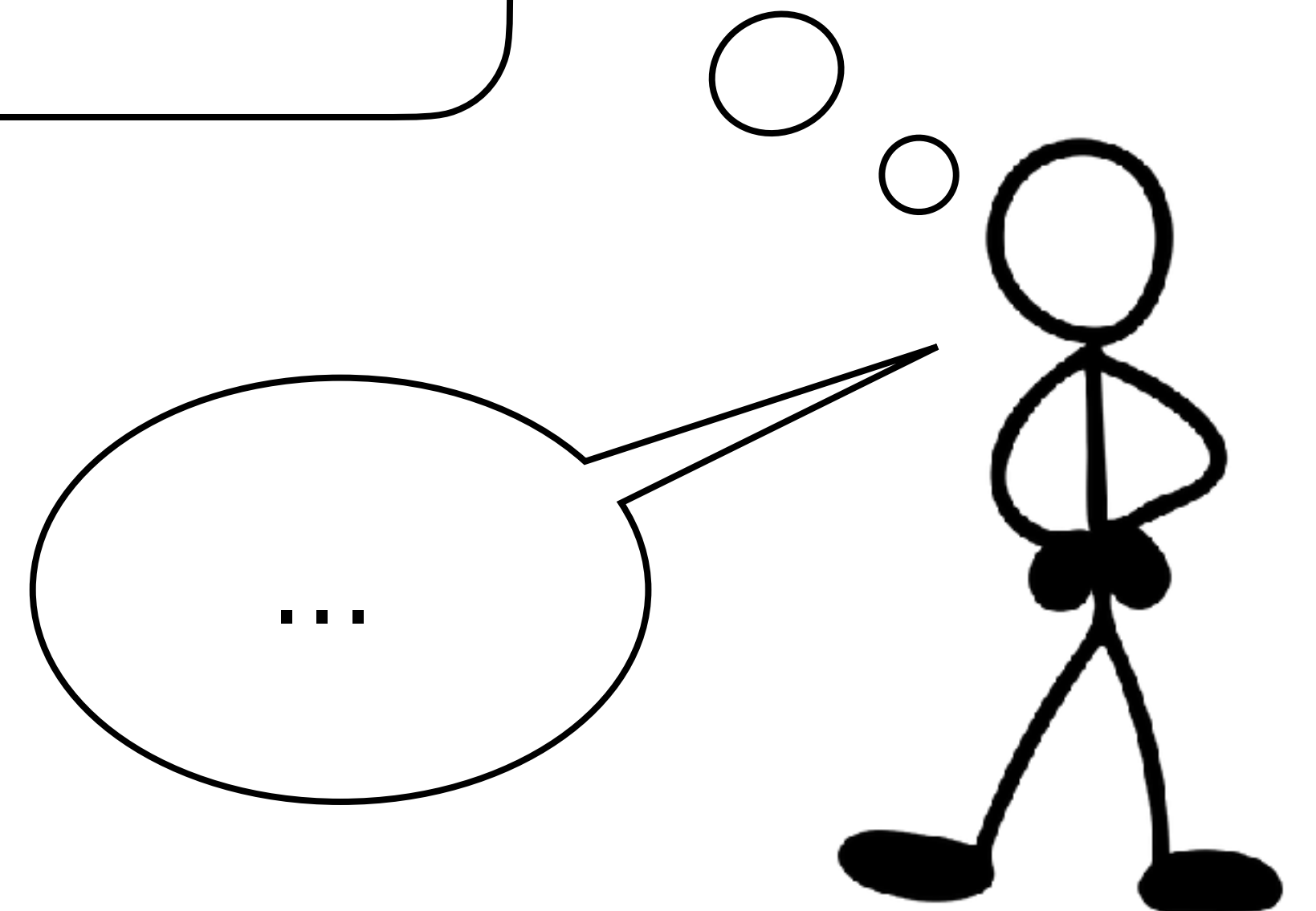
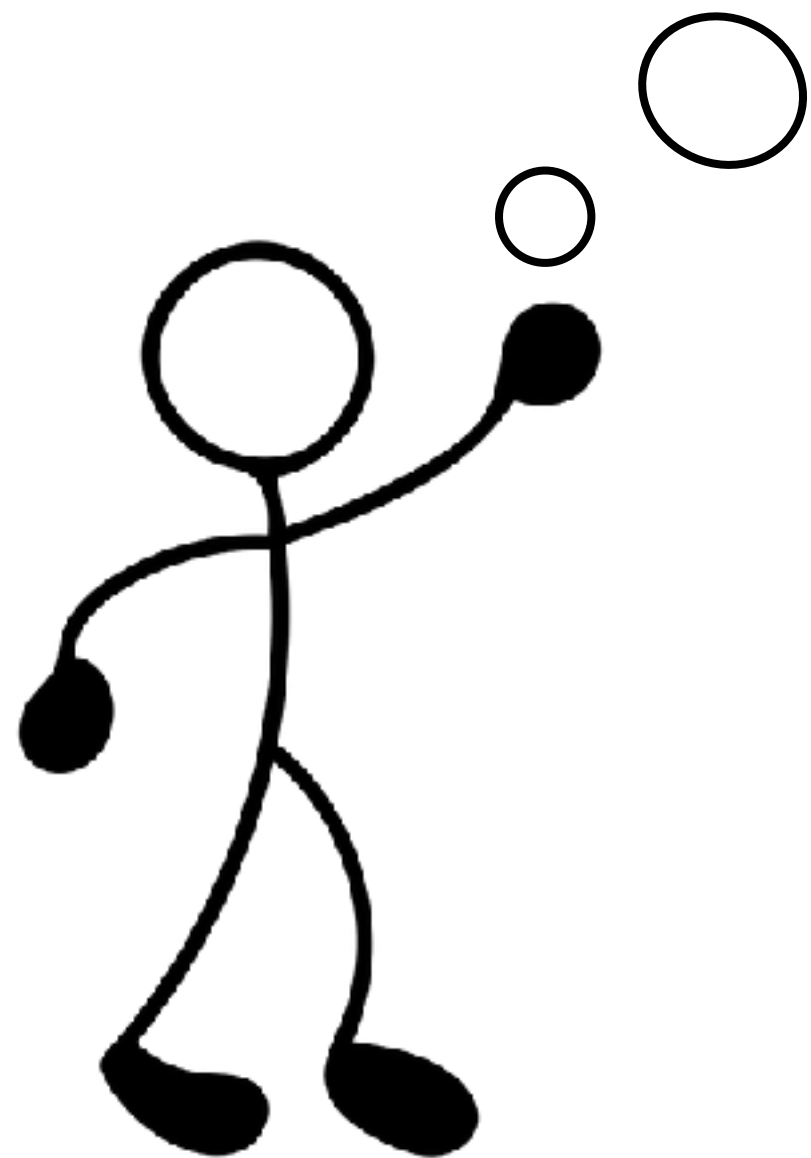
The Gricean Maxims

(Grice, 1975)

A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)



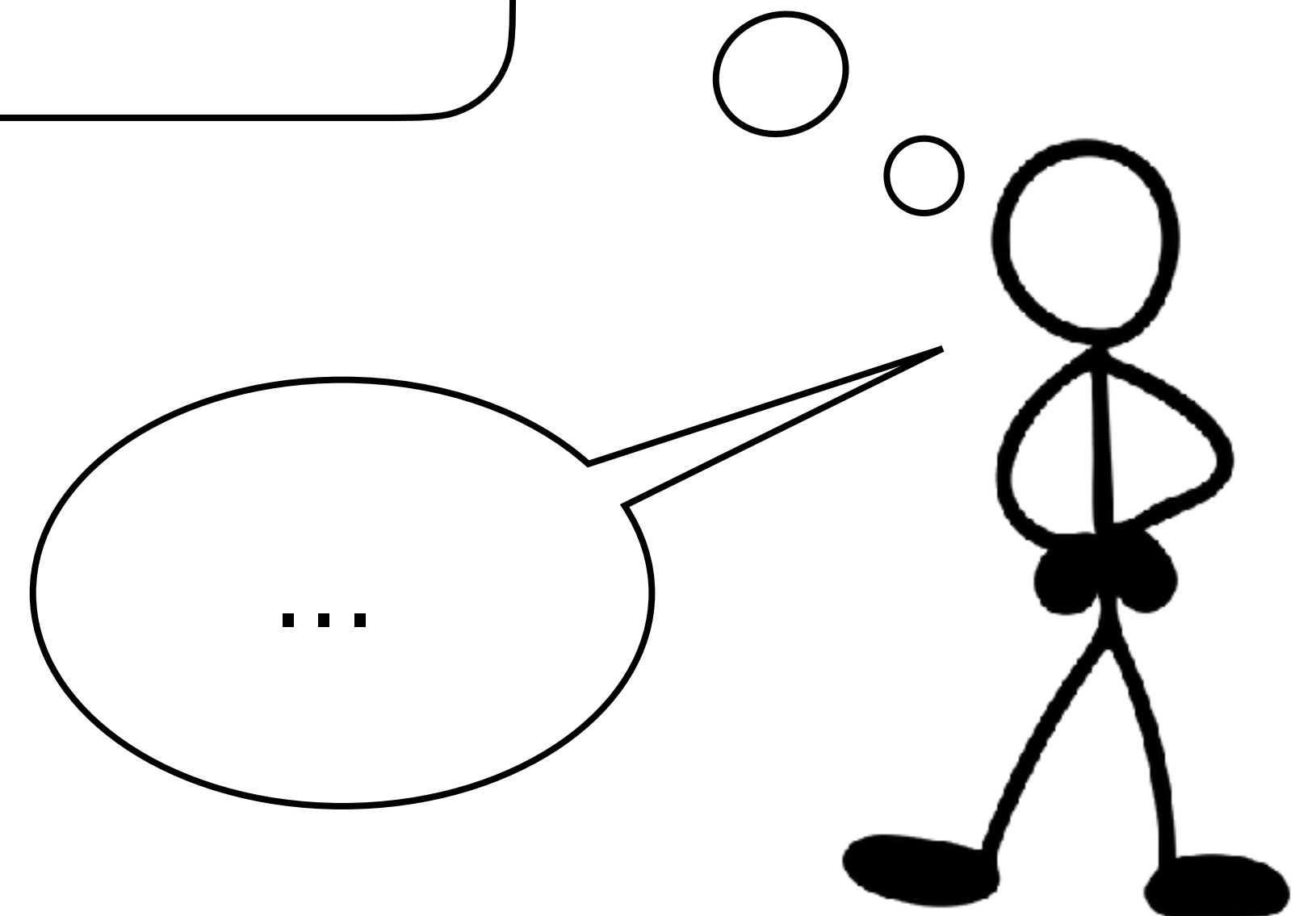
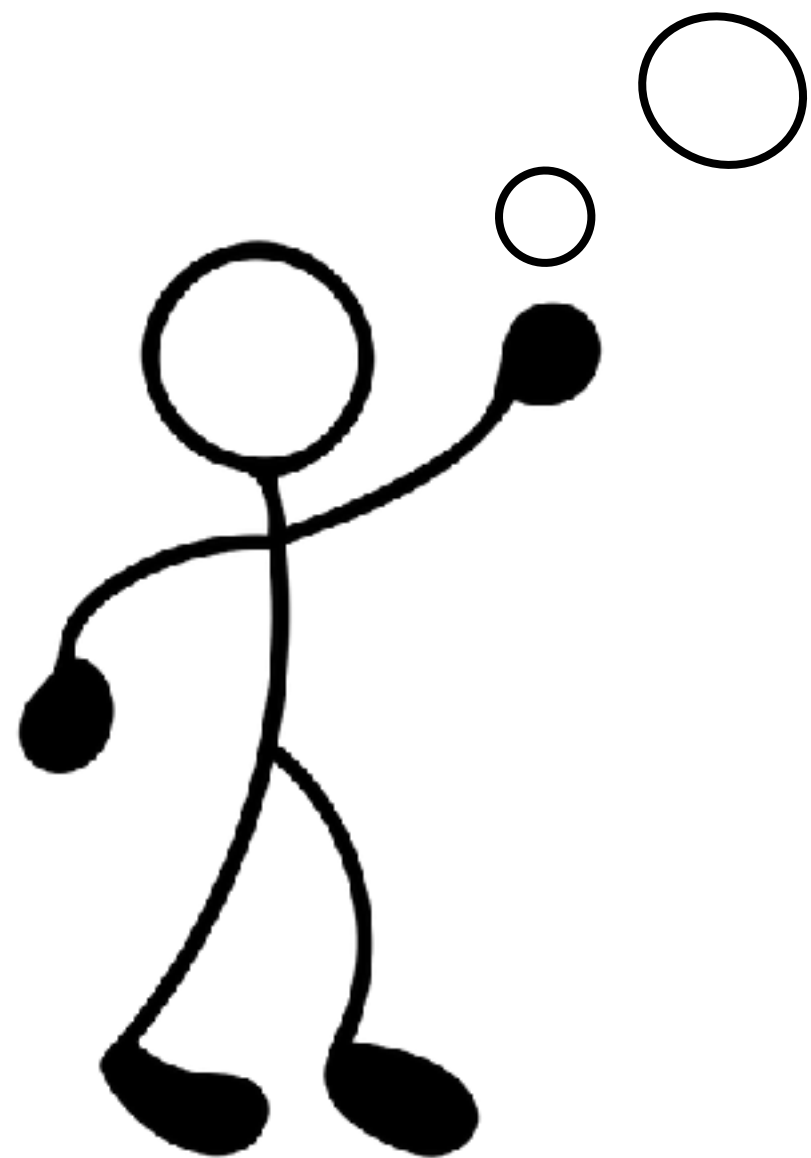
A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)

Maxim of Quantity: Be informative!

- (A) Make your contribution as informative as is required (for the current purposes of the exchange).
- (B) Do not make your contribution more informative than is required.



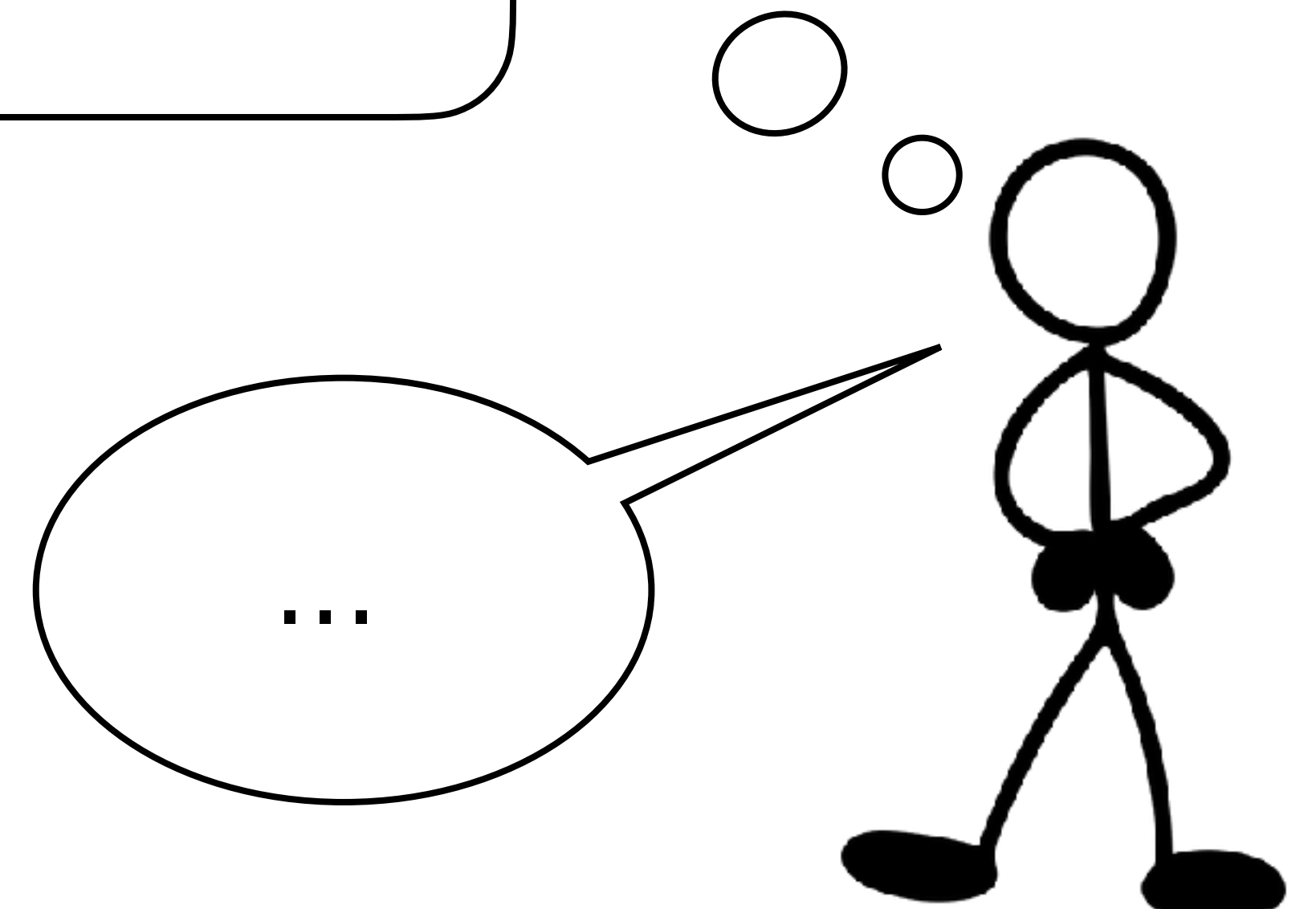
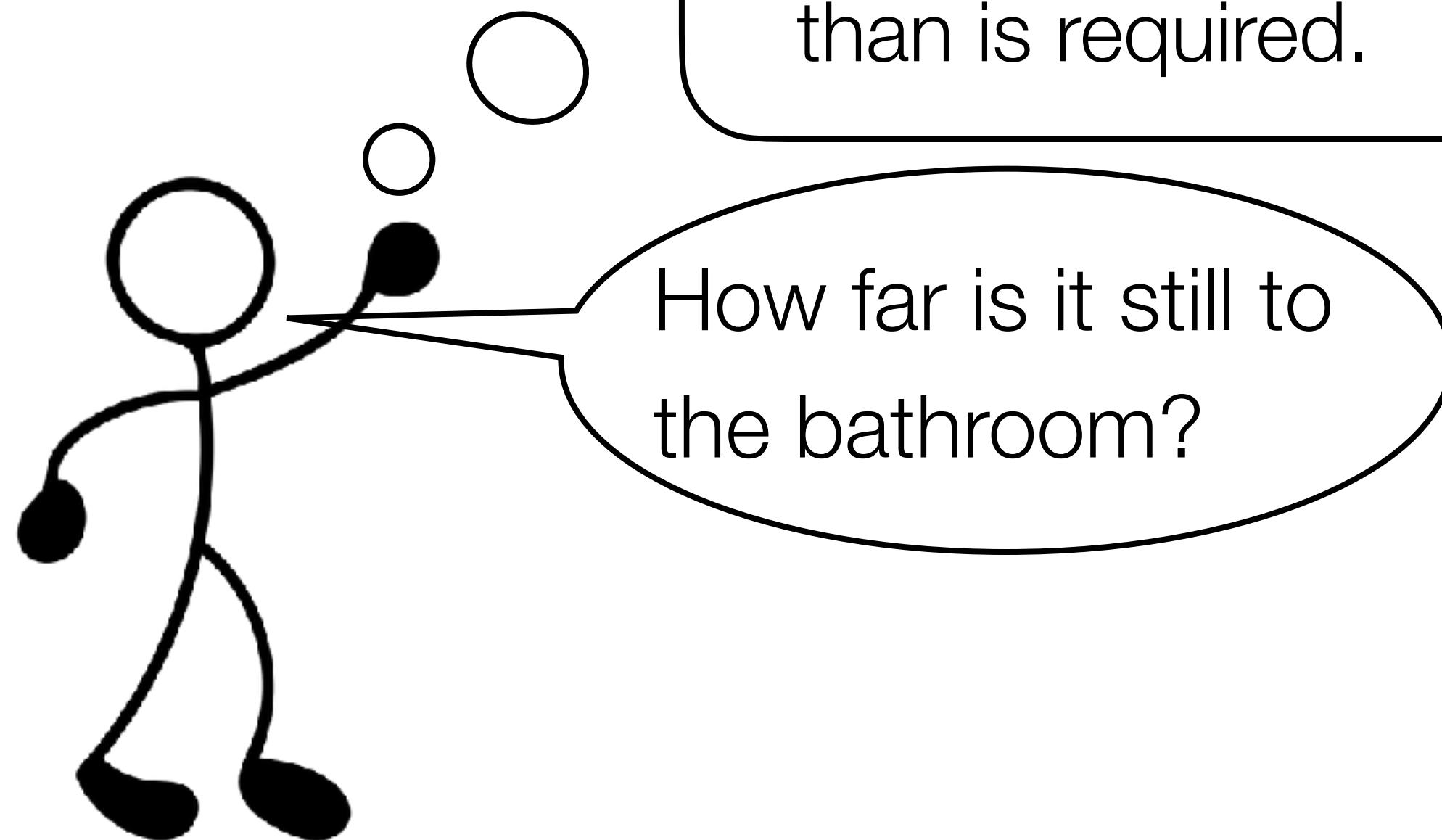
A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)

Maxim of Quantity: Be informative!

- (A) Make your contribution as informative as is required (for the current purposes of the exchange).
- (B) Do not make your contribution more informative than is required.



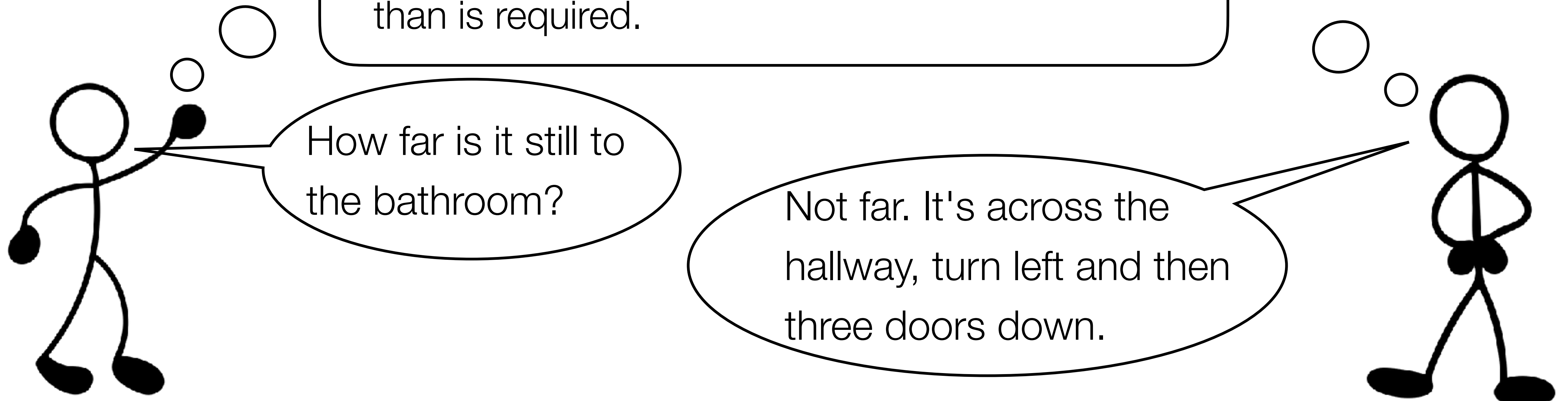
A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)

Maxim of Quantity: Be informative!

- (A) Make your contribution as informative as is required (for the current purposes of the exchange).
- (B) Do not make your contribution more informative than is required.



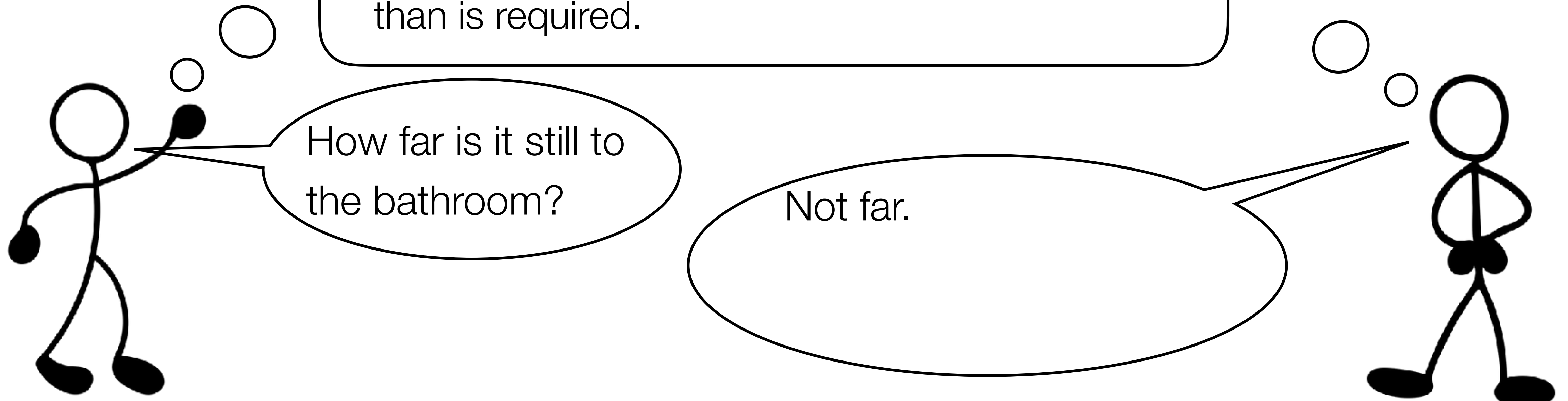
A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)

Maxim of Quantity: Be informative!

- (A) Make your contribution as informative as is required (for the current purposes of the exchange).
- (B) Do not make your contribution more informative than is required.



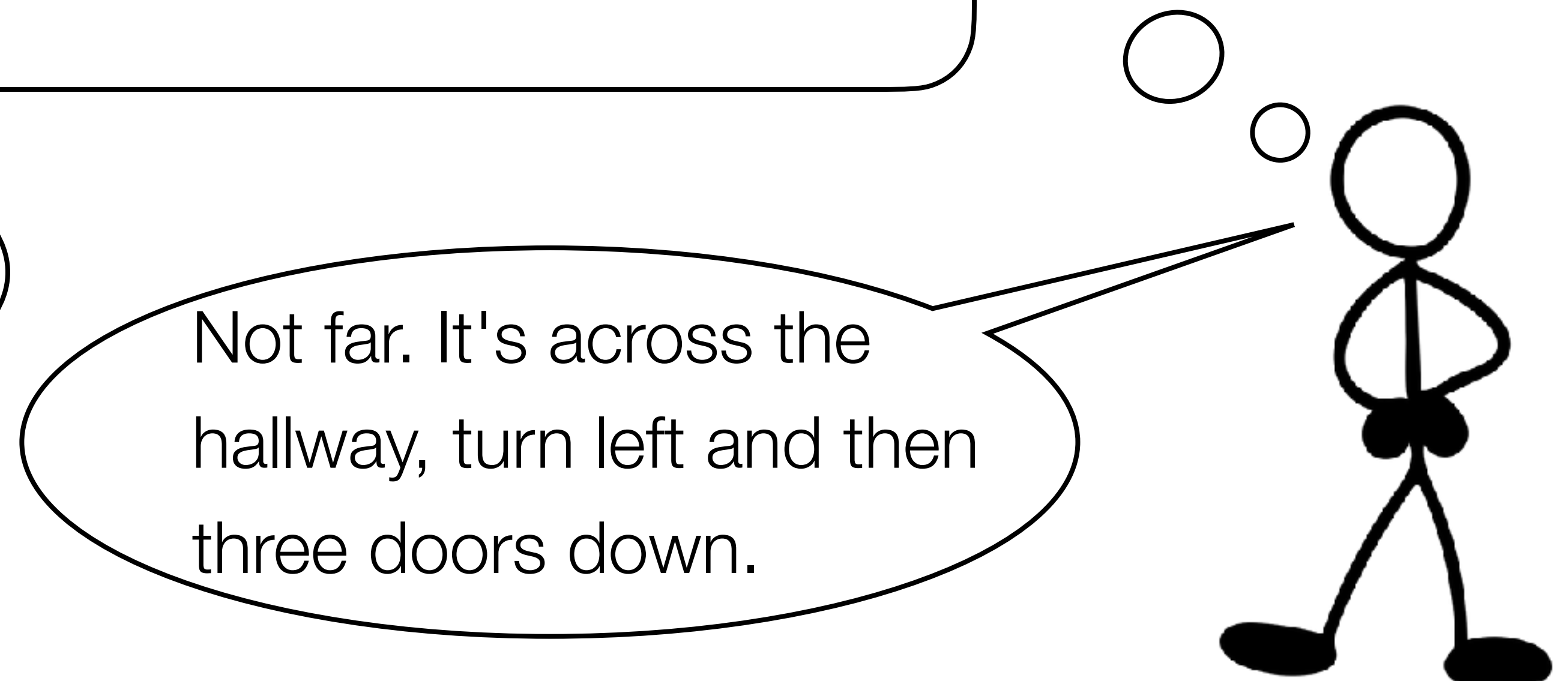
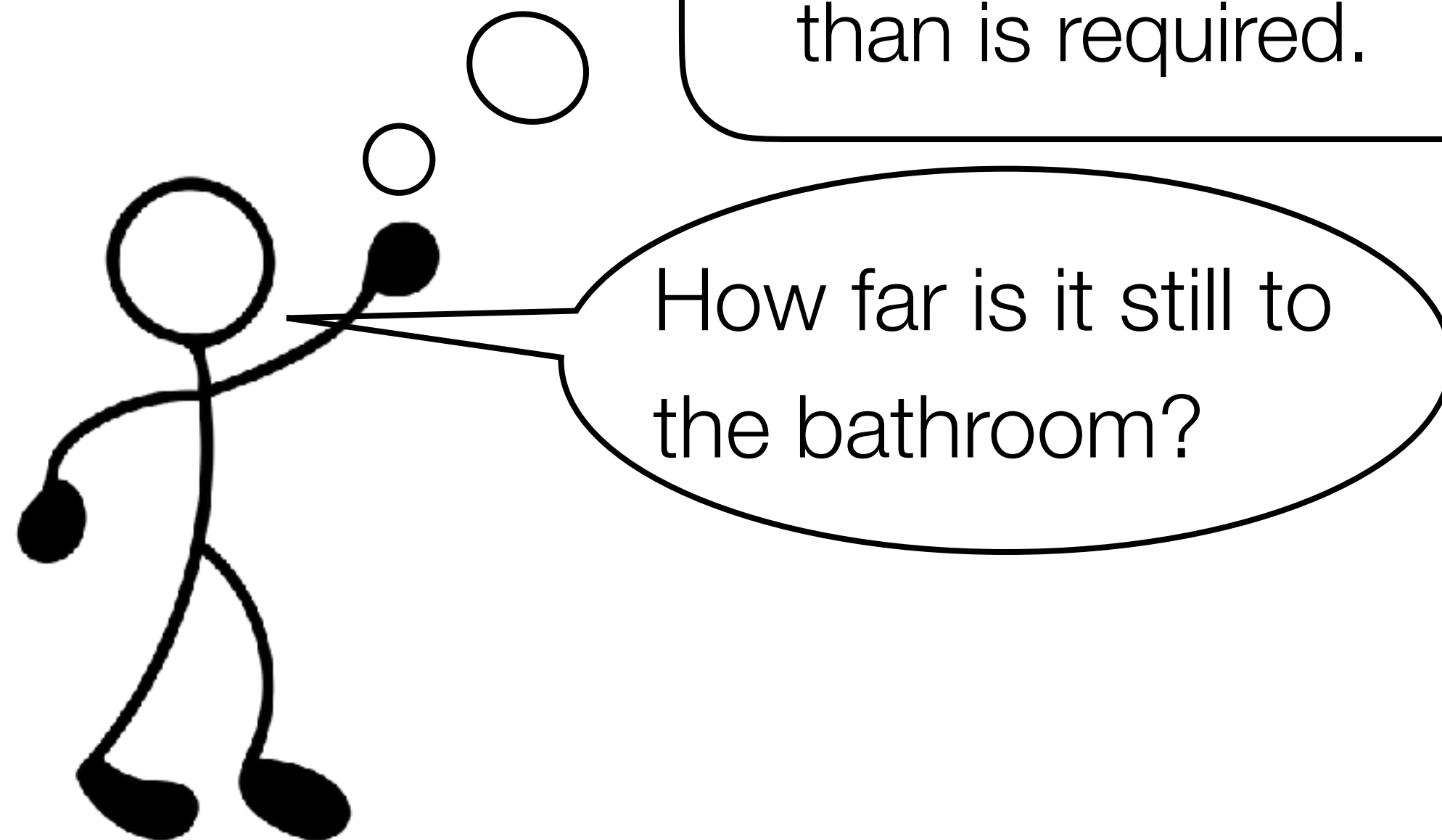
A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)

Maxim of Quantity: Be informative!

- (A) Make your contribution as informative as is required (for the current purposes of the exchange).
- (B) Do not make your contribution more informative than is required.



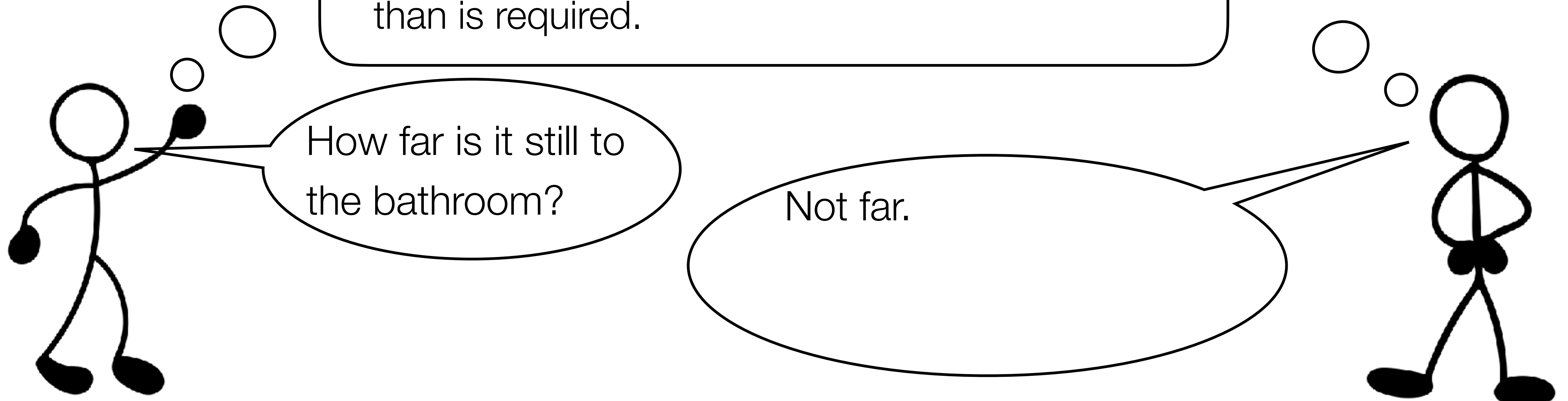
A Pragmatic Approach to Image Descriptions

The Gricean Maxims

(Grice, 1975)

Maxim of Quantity: Be informative!

- (A) Make your contribution as informative as is required (for the current purposes of the exchange).
- (B) Do not make your contribution more informative than is required.



A Pragmatic Approach to Image Descriptions

The Gricean Maxims

Maxim of Quantity: Be informative!

Choosing what is informative depends on **context**.

A Pragmatic Approach to Image Descriptions

The Gricean Maxims

Maxim of Quantity: Be informative!

Choosing what is informative depends on **context**.

There is no one-size-fits-all approach to (accessibility) communication.

[See also: Stangl et al. 2021; Muehlbradt & Kane, 2022; Herskovitz et al. 2023]



WIKIPEDIA
The Free Encyclopedia

The Role of Context for Image Descriptions

Multimodal Pedagogy

Multimodal pedagogy is an approach to the teaching of writing that implements different modes of communication.^{[1][2]} **Multimodality** refers to the use of multiple modes of communication, such as text, images, and sound, in a single communication.

The visual mode conveys meaning through images, diagrams, and other visual elements.

The aural mode refers to sound in the form of speech or music.

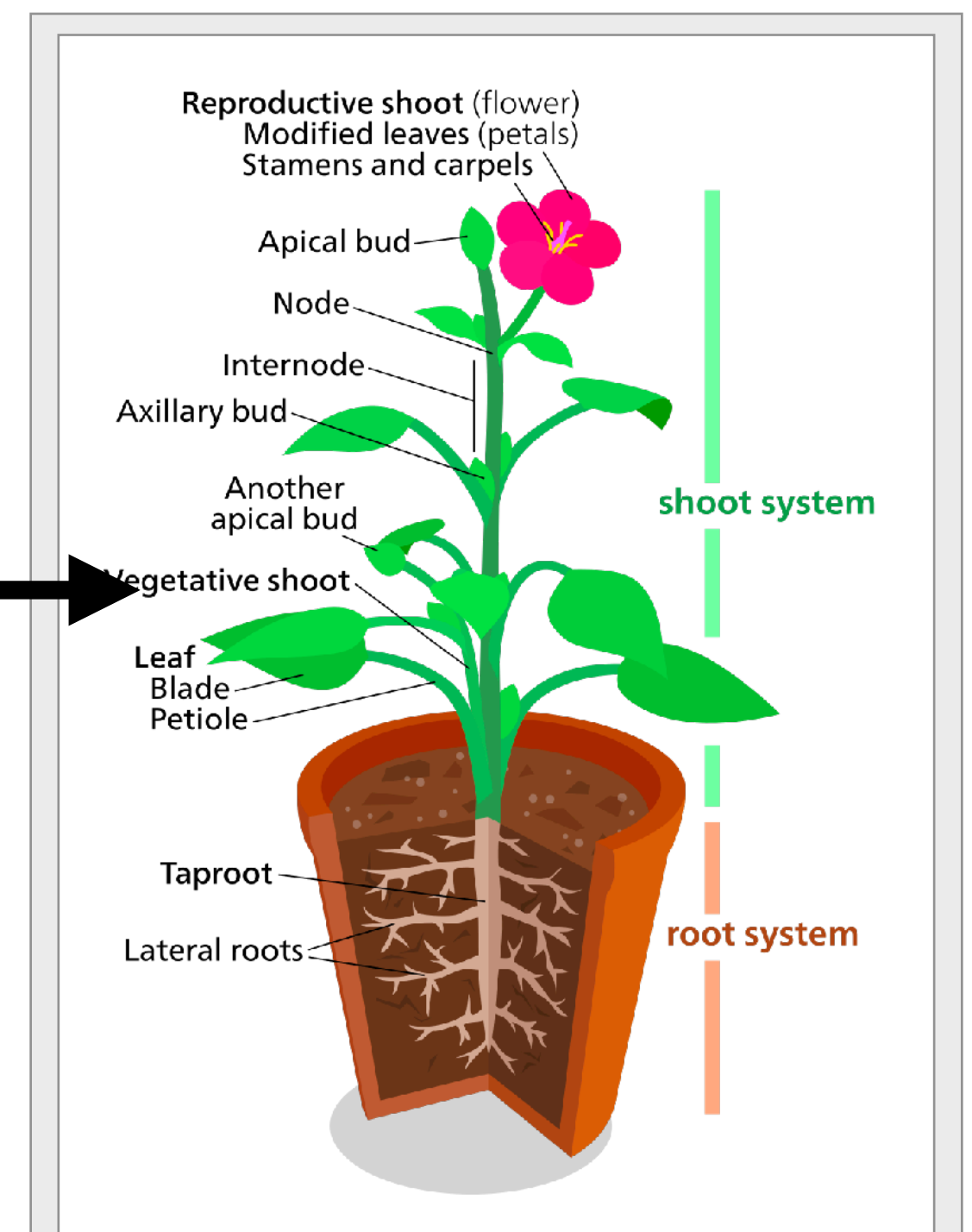
The written mode refers to text in the form of letters, words, and sentences.

The gestural mode refers to physical movement, such as hand gestures or facial expressions.

The multimodal mode refers to the combination of two or more of the above modes.

Multimodality as a term was coined in the 1990s and has since been used as early as **Egyptian hieroglyphs** and classical **rhetoric**.^[7] Compositionists and writing theorists have been exploring how the five modes of communication interact with each other and how multimodality can be used in the teaching of writing since the 20th century.^[8]

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.



Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.



WIKIPEDIA
The Free Encyclopedia

The Role of Context for Image Descriptions

Plant Anatomy

Plant anatomy or **Phytotomy** is the general term for the study of the internal **structure** of **plants**. Originally it included **plant morphology**, the description of the physical form and external structure of plants, but since the mid-20th century plant anatomy has been more concerned with the internal structure of plants. Plant anatomy is now frequently involved in **microscopy**.^[3]

Structural divisions ^[edit]

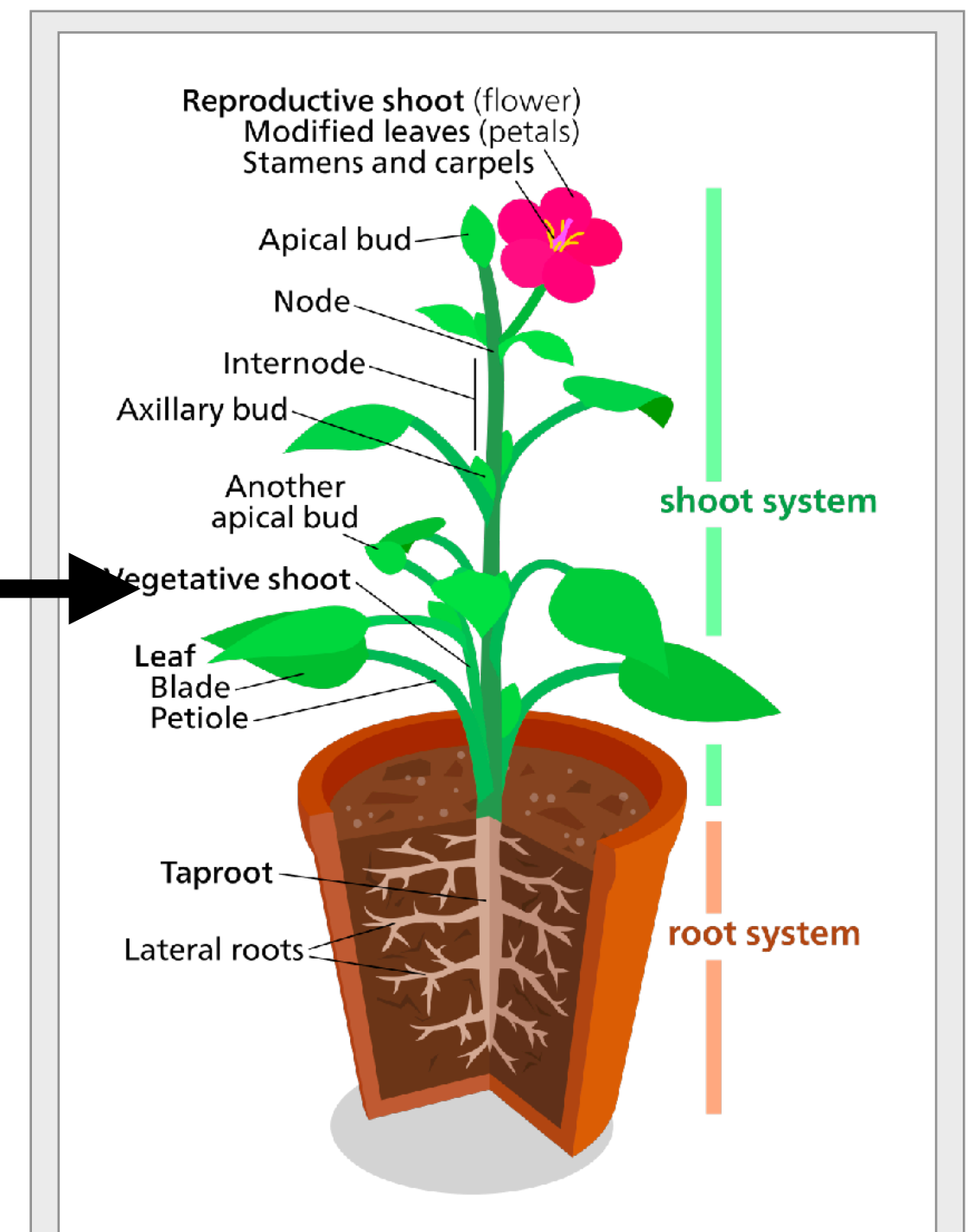
Some studies of plant anatomy use the following structural divisions: nutrient transport, flowering, pollination, and growth. Plant anatomy is divided into the following structural categories:

Flower anatomy, including study of the **Calyx**, **Corolla**, **Androecium**, and **Gynoecium**

Leaf anatomy, including study of the **Epidermis**, **stomata** and **Palisade cells**

Stem anatomy, including **Stem structure** and **vascular tissues**, **buds** and **shoot apex**

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.



A diagram of the anatomy of a plant with labels of structural parts of the plant and the roots.

Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Kreiss, Bennett, Hooshmand, Zelikman, Morris, Potts (EMNLP 2022)



WIKIPEDIA
The Free Encyclopedia

The Role of Context for Image Descriptions

Plant Anatomy

Plant anatomy or **Phytotomy** is the general term for the study of the internal **structure** of **plants**. Originally it included **plant morphology**, the description of the physical form and external structure of plants, but since the mid-20th century plant anatomy has been more concerned with the internal structure. Plant anatomy is now frequently involved in **microscopy**.^[3]

Structural divisions ^[edit]

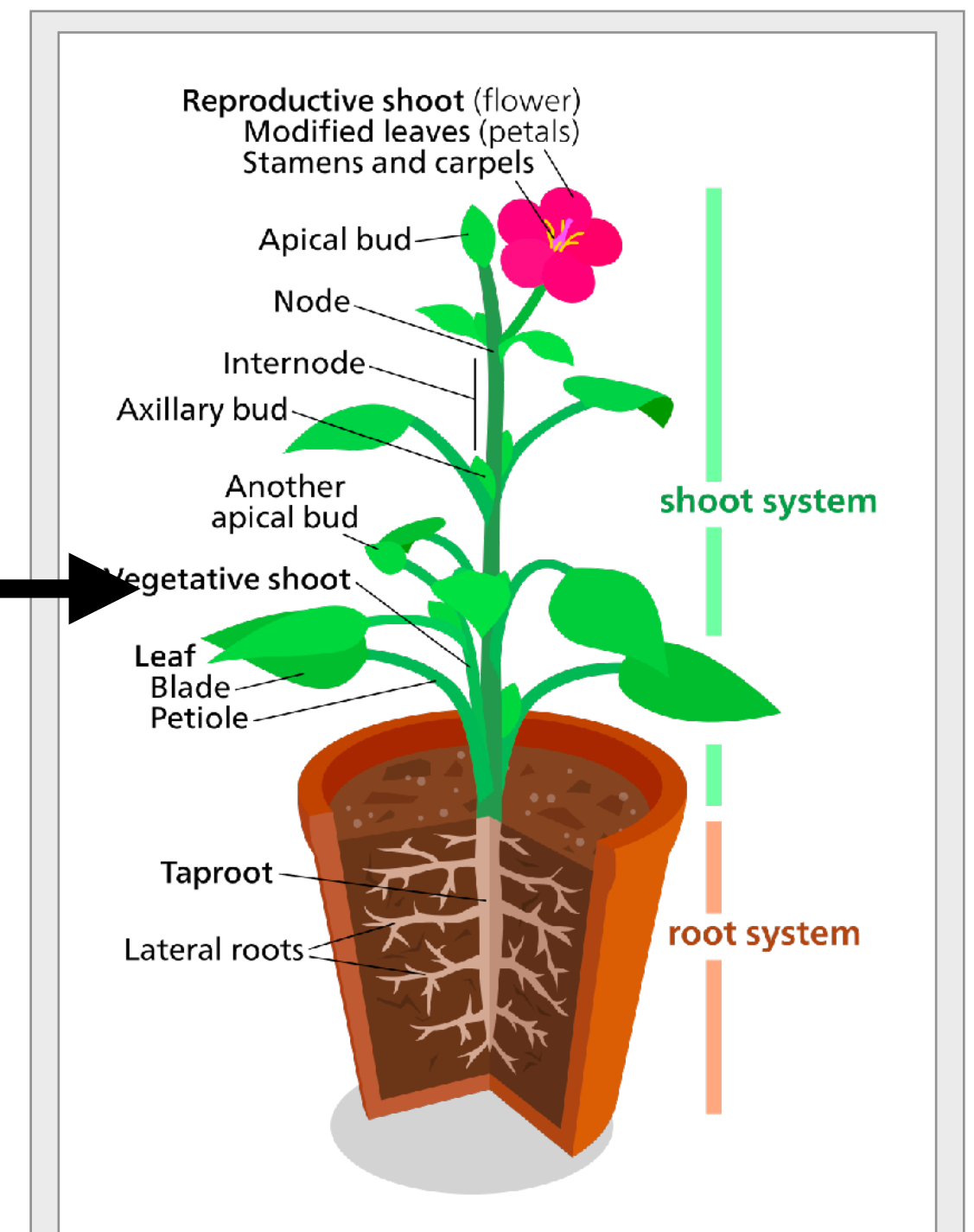
Some studies of plant anatomy use the following structural divisions: nutrient transport, flowering, pollination, and growth. They are divided into the following structural categories:

Flower anatomy, including study of the **Calyx**, **Corolla**, **Androecium**, and **Gynoecium**

Leaf anatomy, including study of the **Epidermis**, **stomata** and **Palisade cells**

Stem anatomy, including **Stem structure** and **vascular tissues**, **buds** and **shoot apex**

A sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. Part of the root system are the taproots and lateral roots. The taproot refers to the central root and the lateral roots are the smaller side roots that ...



A diagram of the anatomy of a plant with labels of structural parts of the plant and the roots.

Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Kreiss, Bennett, Hooshmand, Zelikman, Morris, Potts (EMNLP 2022)

The **Image**'s Communicative Goal, or: **Context Matters!**

Data

Modeling

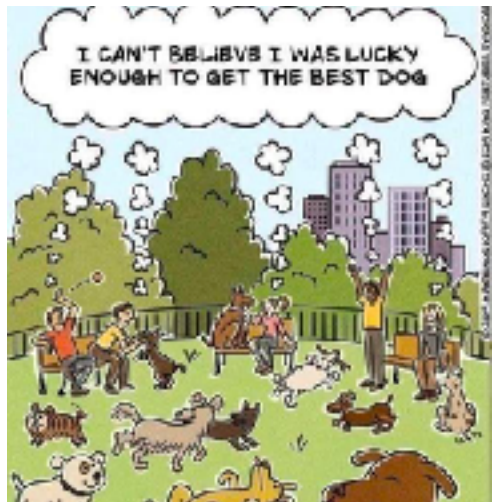
Evaluation

The **Image**'s Communicative Goal, or: **Context Matters!**

Data



a dog
sitting



cartoon
of a dog
park



a
dinosaur
telling a
joke

e.g., MSCOCO (Lin et al., 2014); Flickr30k (Hodosh et al. 2013); VizWizCaptions (Gurari et al. 2020)

Modeling

Evaluation

The **Image**'s Communicative Goal, or: **Context Matters!**

Data



image-
based
text

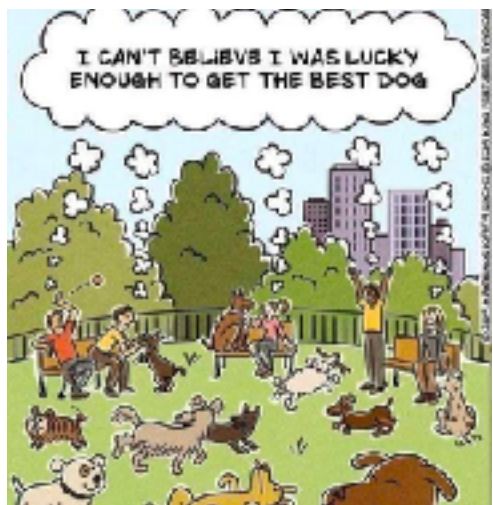


image-
based
text



image-
based
text

e.g., MSCOCO (Lin et al., 2014); Flickr30k (Hodosh et al. 2013); VizWizCaptions (Gurari et al. 2020)

Modeling



Text
Generation
Model



image-based
text

general principle in, e.g., Xu et al. 2015, Li et al. 2020, Hossain et al. 2021, Radford et al. 2021

Evaluation

The **Image**'s Communicative Goal, or: **Context Matters!**

Data



image-
based
text

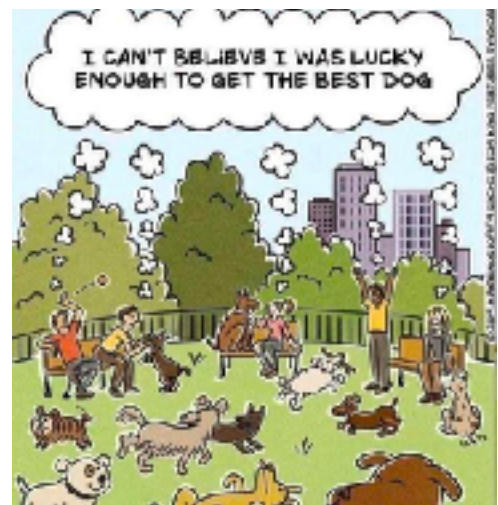


image-
based
text



image-
based
text

e.g., MSCOCO (Lin et al., 2014); Flickr30k (Hodosh et al. 2013); VizWizCaptions (Gurari et al. 2020)

Modeling



Text
Generation
Model

image-based
text

general principle in, e.g., Xu et al. 2015, Li et al. 2020, Hossain et al. 2021, Radford et al. 2021

Evaluation

image-
based
text



Scoring Model

score

Hessel et al. 2021, Lee et al. 2021a, Feinglass & Yang 2021, Lee et al. 2021b

The **Image**'s Communicative Goal, or: **Context Matters!**

Data

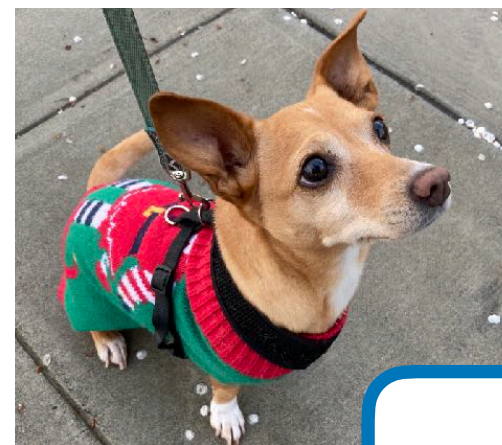


image-based
text

+ context

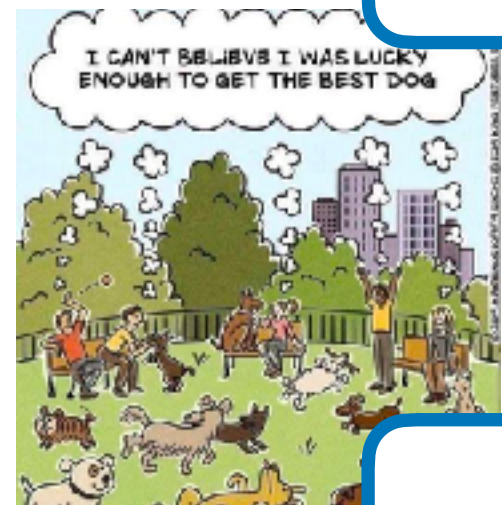


image-based
text

+ context

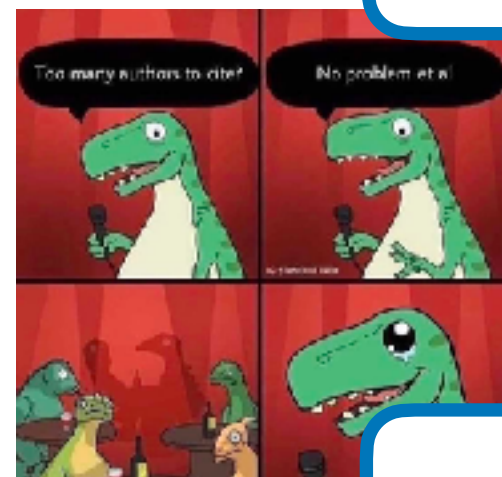


image-based
text

+ context

e.g., MSCOCO (Lin et al., 2014); Flickr30k (Hodosh et al. 2013); VizWizCaptions (Gurari et al. 2020)

context



Modeling

Text
Generation
Model

image-based
text

general principle in, e.g., Xu et al. 2015, Li et al. 2020, Hossain et al. 2021, Radford et al. 2021

context

image-based
text



Evaluation

Scoring Model

score

Hessel et al. 2021, Lee et al. 2021a, Feinglass & Yang 2021, Lee et al. 2021b

The **Image**'s Communicative Goal, or: **Context Matters!**



Concadia: A naturalistic dataset containing rich contextual information.

96,918 images with captions, alt descriptions and surrounding paragraph from Wikipedia

Contextual information: (1) domain, (2) article title, (3) closest paragraph, (4) description / caption

Wikipedia-Article on **Banana**

image context: In global commerce in 2009, by far the most important cultivars belonged to the triploid AAA group of *Musa acuminata*, commonly referred to as Cavendish group bananas. They accounted for the majority of banana exports, despite only coming into existence in 1836. The cultivars Dwarf Cavendish and Grand Nain (Chiquita Banana) gained popularity in the 1950s after the previous mass-produced cultivar, Gros Michel (also an AAA ...



(Accessibility) **Description:**

Grocery store photo of several bunches of bananas

(Contextualizing) **Caption:**

Cavendish bananas are the main commercial banana cultivars sold in the world market.

Making Models Context-Sensitive

ResNet + LSTM

Input

Encoding

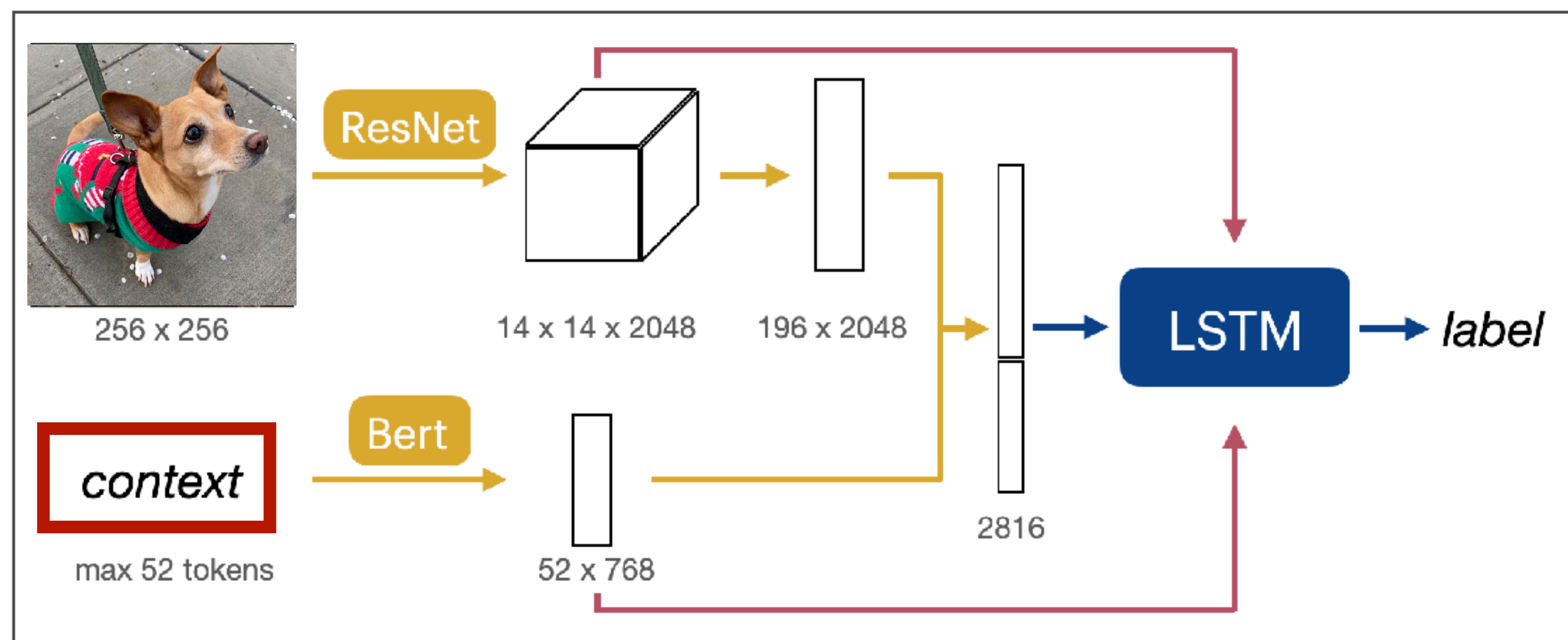
DenseNet + LSTM

Attention

Decoding

Output

OSCAR (VinVL)



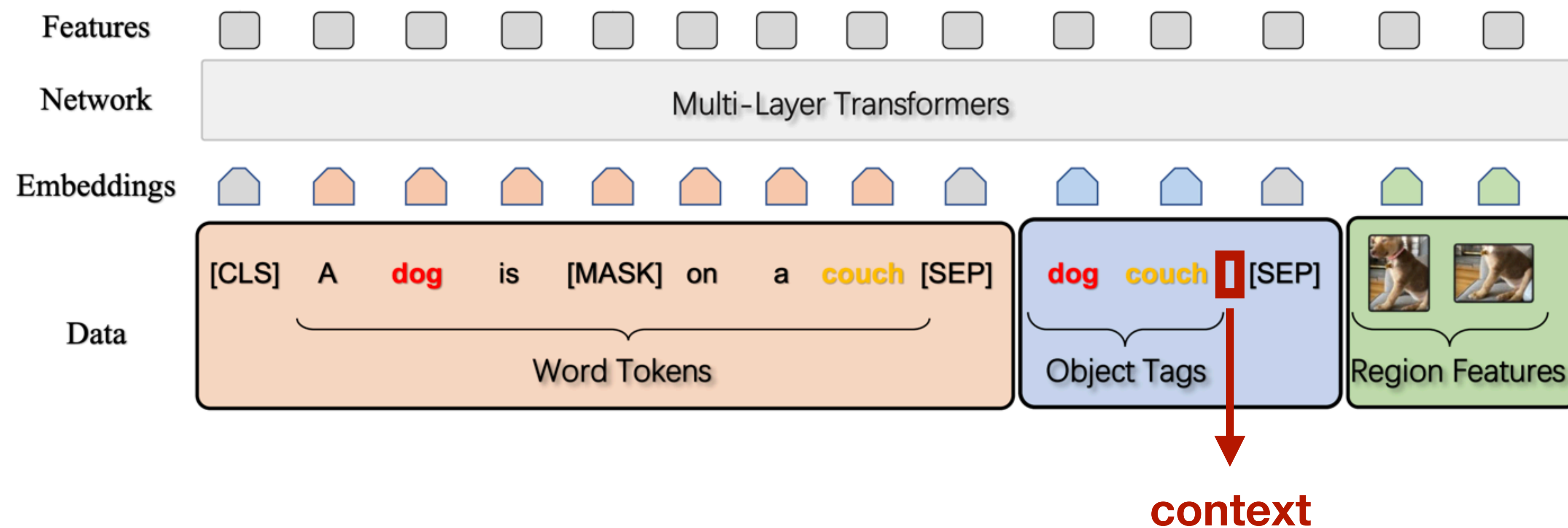
Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

Making Models Context-Sensitive

ResNet + LSTM

DenseNet + LSTM

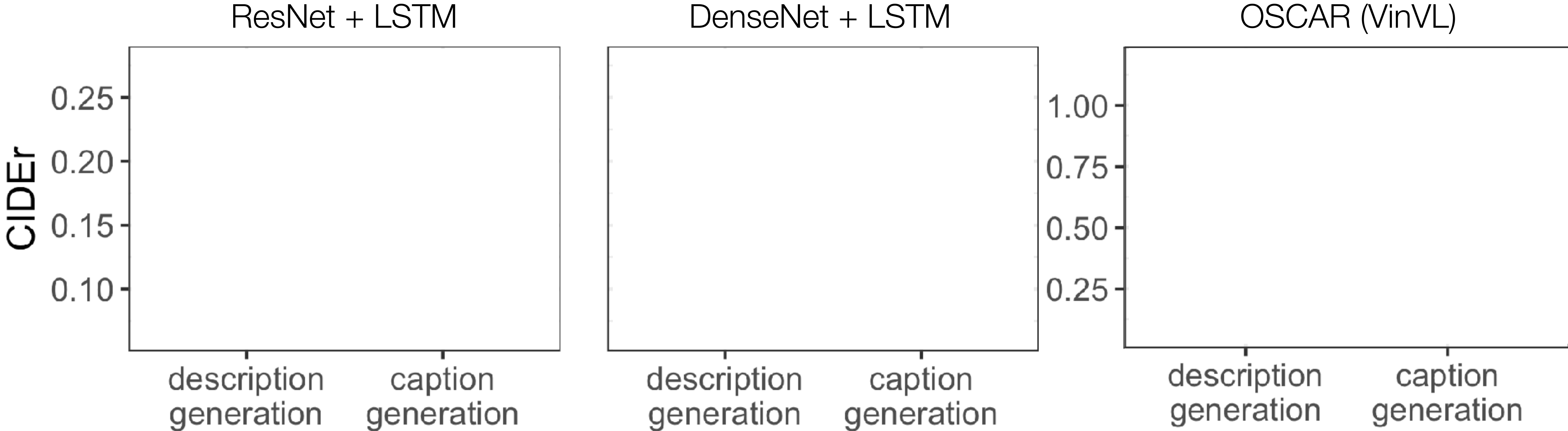
OSCAR (VinVL)



Concadia: Towards image-based text generation with a purpose
Kreiss, Fang, Goodman, Potts (EMNLP 2022)

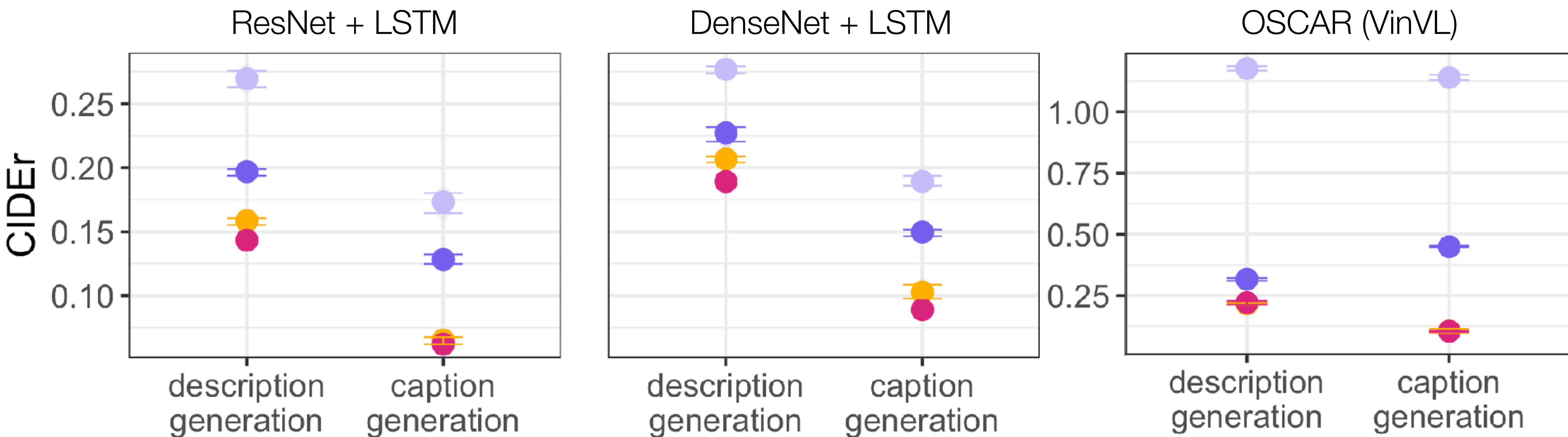
Context Improves Image-Based Text Generation

context none randomized paragraph paragraph caption / description



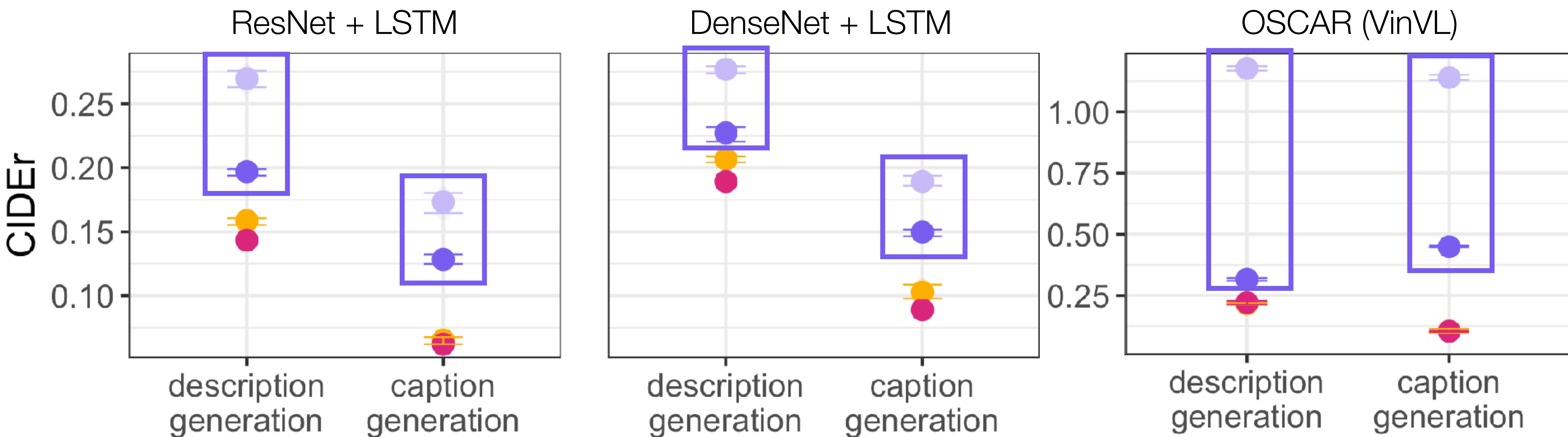
Context Improves Image-Based Text Generation

context ● none ● randomized paragraph ● paragraph ● caption / description



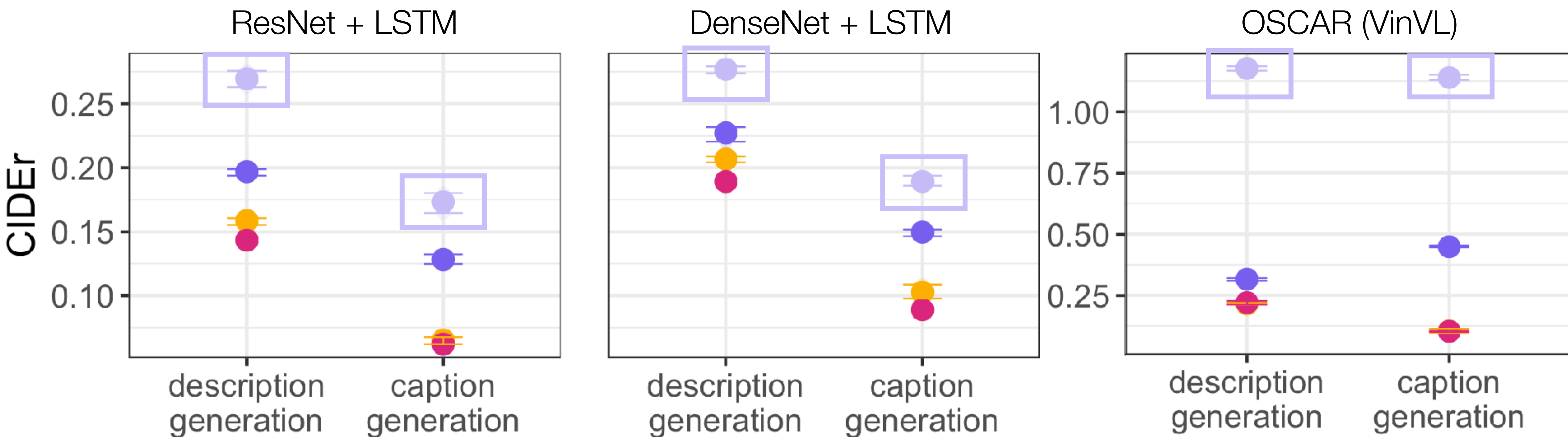
Context Improves Image-Based Text Generation

context ● none ● randomized paragraph ● paragraph ● caption / description



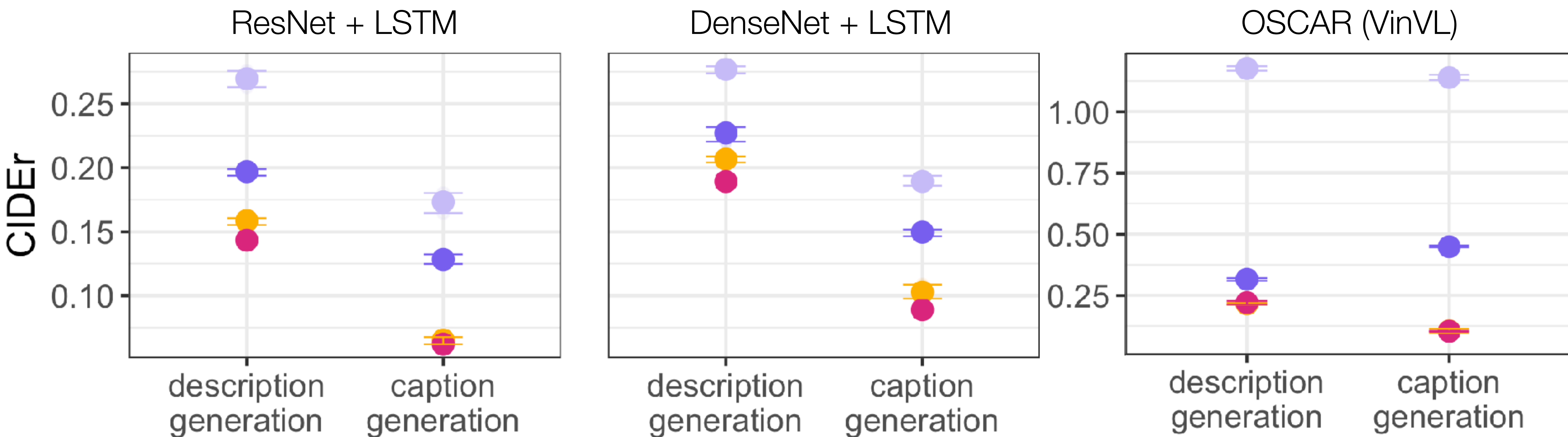
Context Improves Image-Based Text Generation

context none randomized paragraph paragraph caption / description



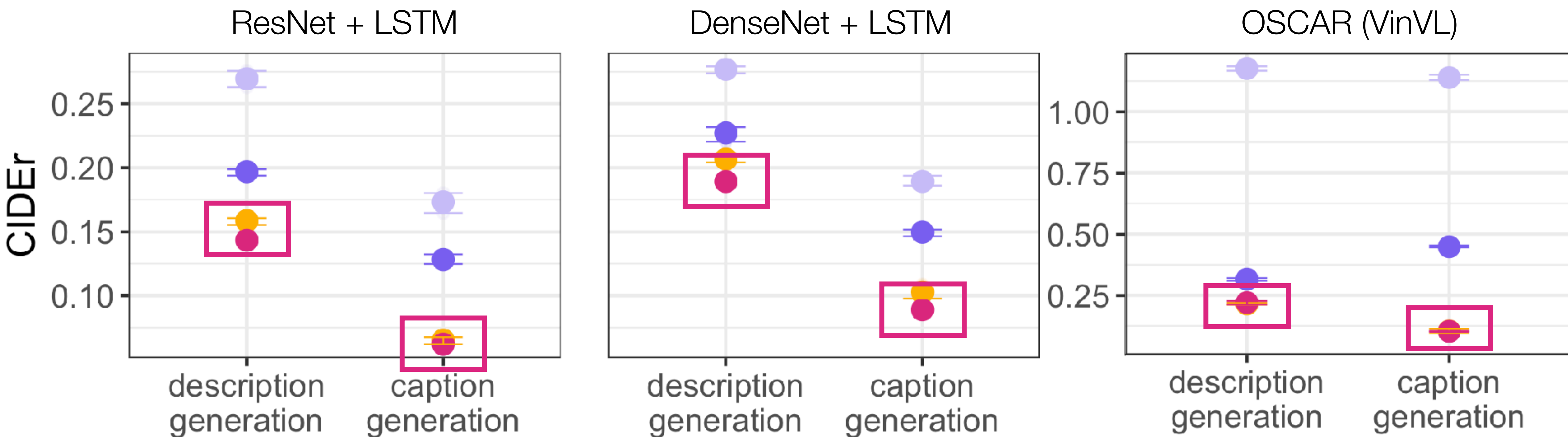
Context Improves Image-Based Text Generation

context ● none ● randomized paragraph ● paragraph ● caption / description



Context Improves Image-Based Text Generation

context ● none ● randomized paragraph ● paragraph ● caption / description



The **Image**'s Communicative Goal, or: **Context Matters!**



Concadia: A naturalistic dataset containing rich contextual information.



Models of image-based text generation:

Supplying contextual information benefits model performance.

Model Evaluation

Which model is better?



Neural
Network
Model (1)



some
image-based text
output



Neural
Network
Model (2)



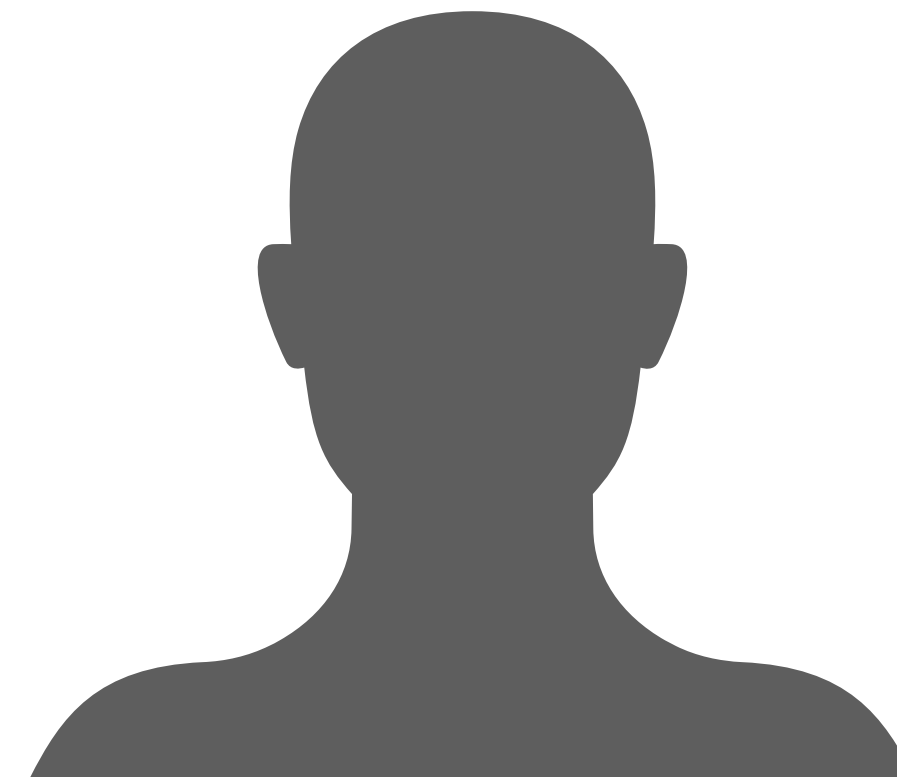
some other!
image-based text
output

Model Evaluation Methods: People

Which model is better?

Neural
Network
Model (1)

Neural
Network
Model (2)



Model Evaluation Methods: People

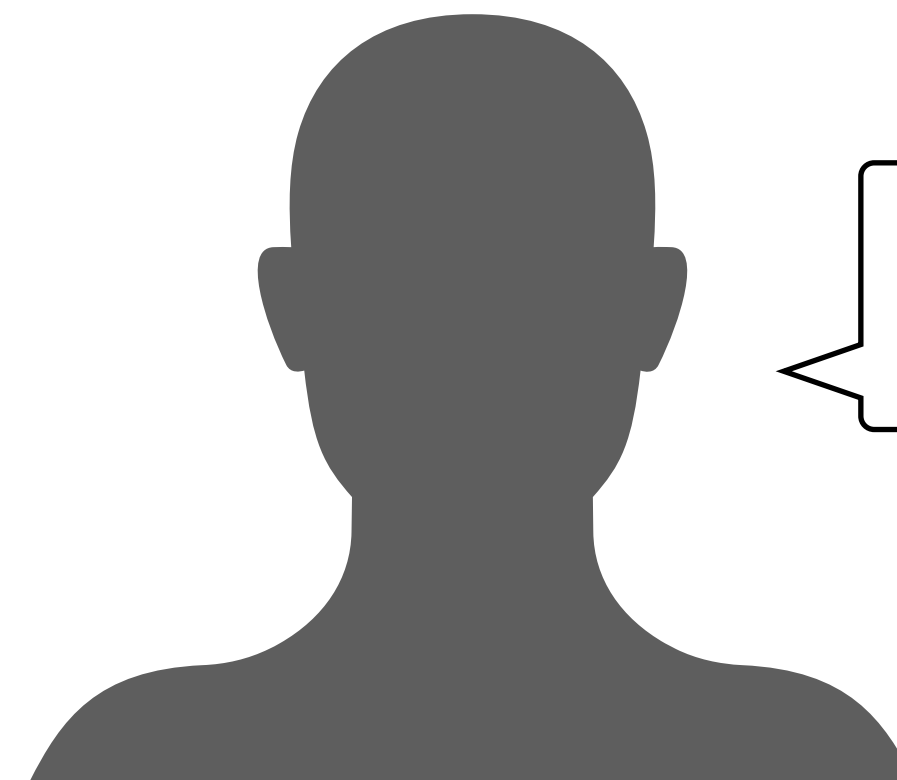
Which model is better?

Neural
Network
Model (1)

Neural
Network
Model (2)



some
image-based text
output



Great!

Model Evaluation Methods: People

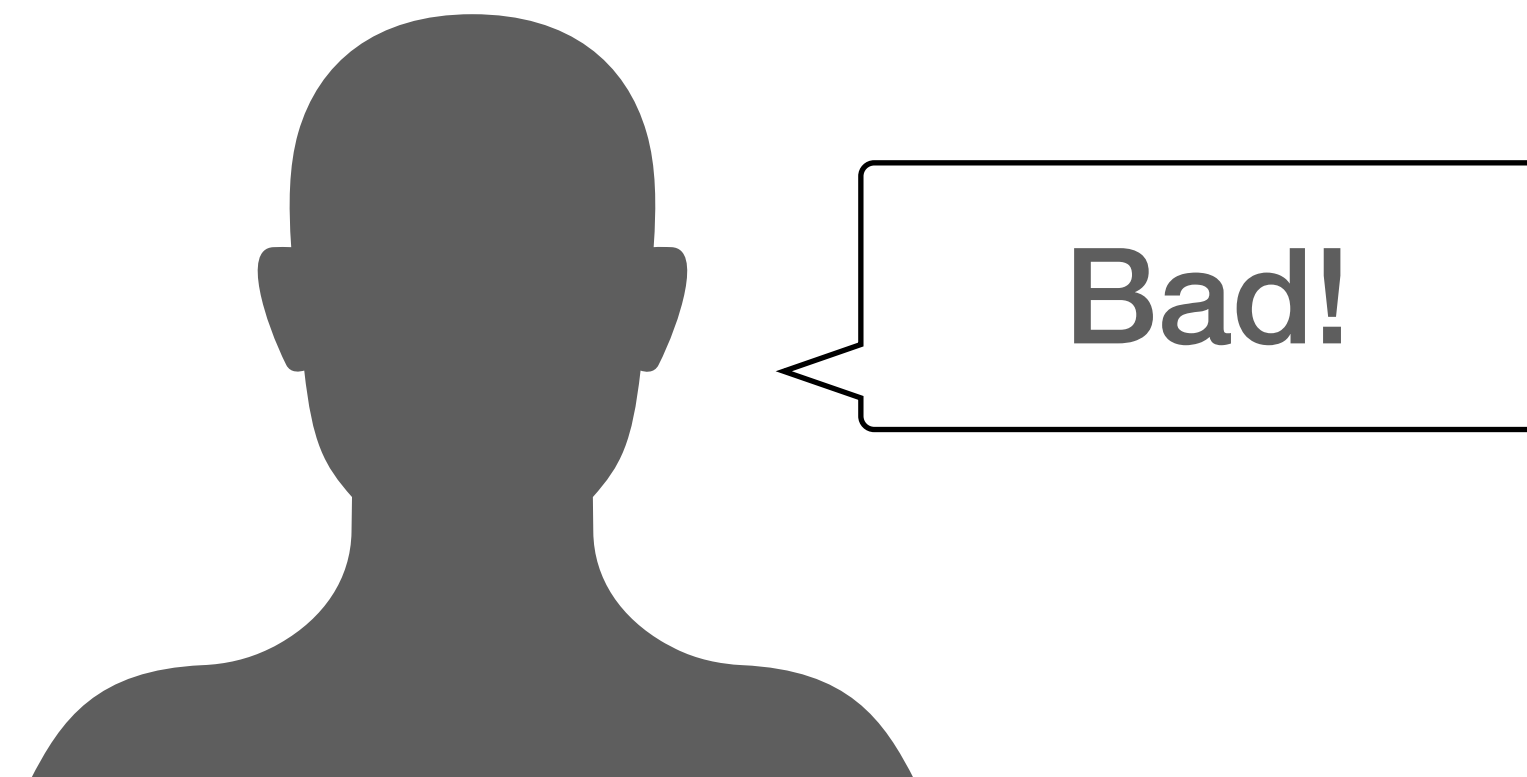
Which model is better?

Neural
Network
Model (1)

Neural
Network
Model (2)



some other!
image-based text
output

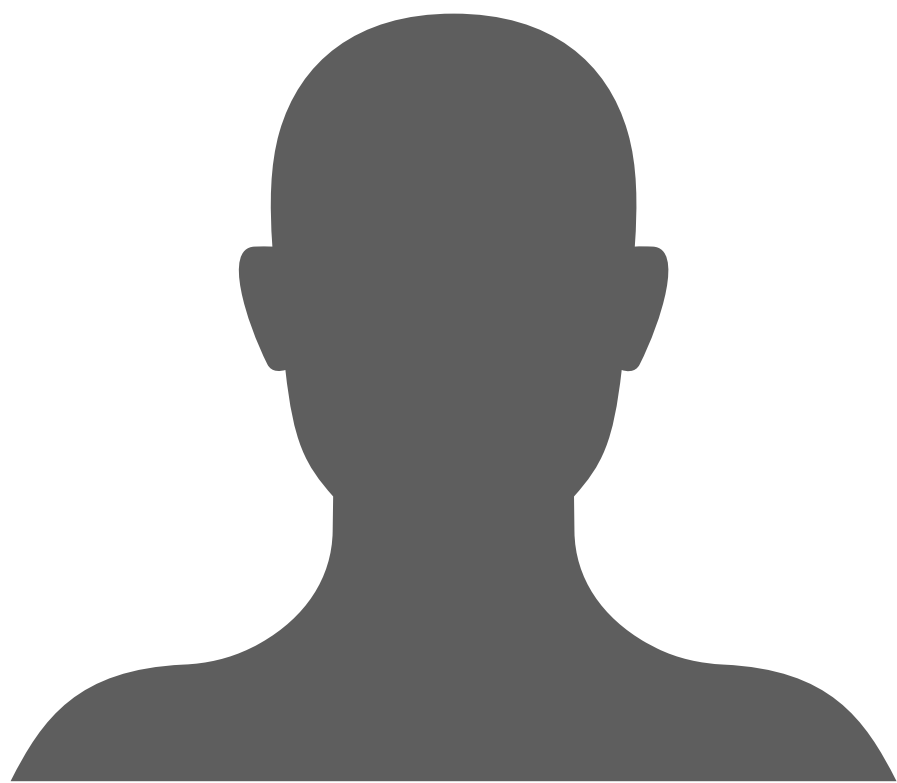


Model Evaluation Methods: People

Which model is better?

Neural
Network
Model (1)

Neural
Network
Model (2)



Pro: Effective!

Con: Slow and costly → not feasible for model development

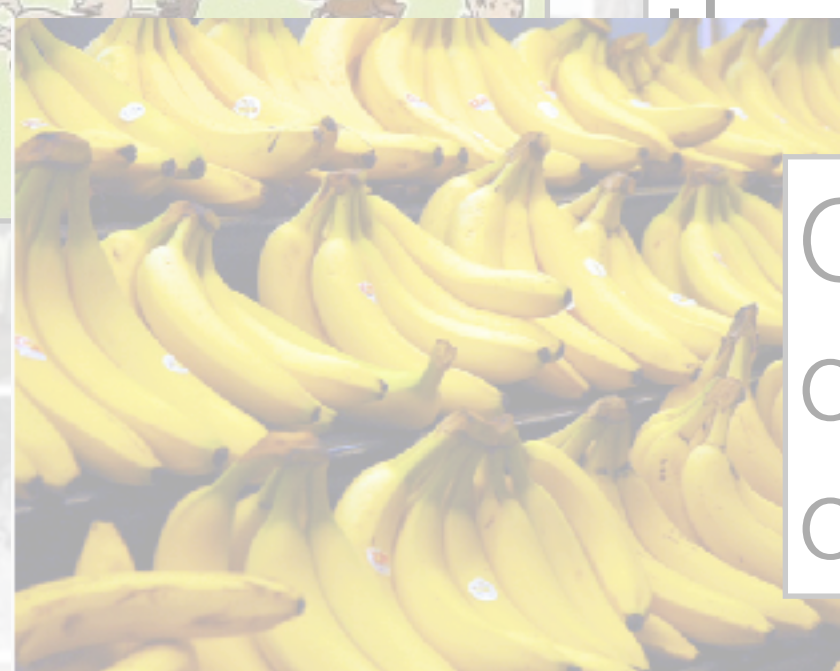
Model Evaluation Methods: Automatic



small dog sitting
outside

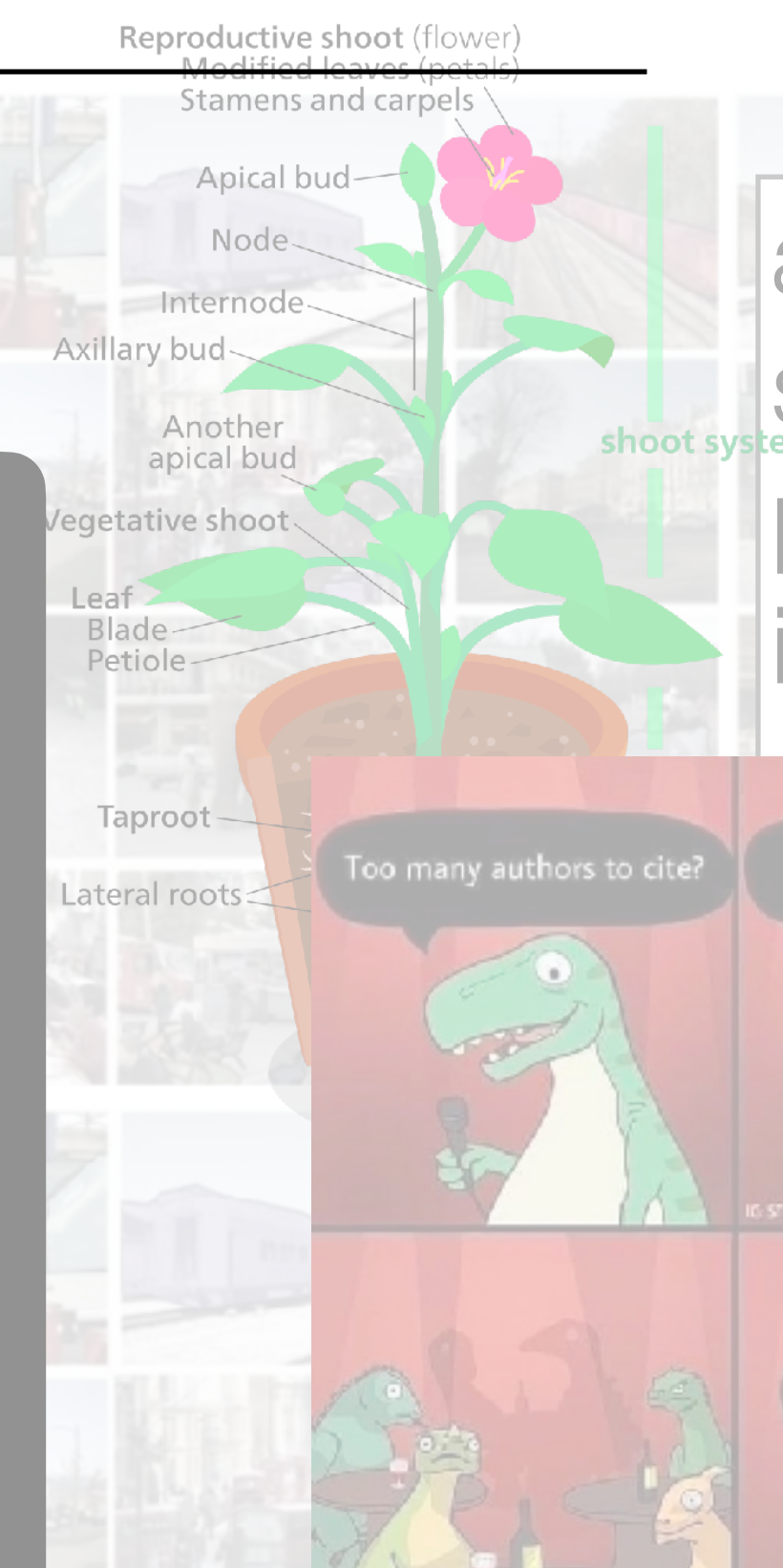


Cartoon of various dogs
and people in a park and
all people thinking the
same thing: I can't believe
I was lucky enough to get
the best dog.



Grocery store photo
of several bunches
of bananas

Model



an educational
sketch of a
plant's anatomy,
including the
shoot system and root

a
dinosaur
telling a
joke

Model Evaluation Methods: Automatic



some image-
based text



Evaluation
Model



Score

Model Evaluation Methods: Automatic

State-of-the-art metric (e.g., CLIPScore)

Hessel et al. 2021, Kasai et al. 2021



Evaluation
Model



Score

some image-
based text



Model Evaluation Methods: Automatic

State-of-the-art metric (e.g., CLIPScore): Inherently context-**in**dependent



some image-
based text



Evaluation
Model



Score

The Role of Context for Evaluation



Context matters for image descriptions for accessibility: Challenges for referenceless evaluation metrics

115 Kreiss, Bennett, Hooshmand, Zelikman, Morris, Potts (EMNLP 2022)

The Role of Context for Evaluation

Building Material



Building material is material used for construction. Many naturally occurring substances, such as clay, rocks, sand and wood, even twigs and leaves, have been used to construct buildings. Apart from naturally occurring materials, ...

Christian Cross



The Christian cross, seen as a representation of the instrument of the crucifixion of Jesus, is the best-known symbol of Christianity. It is related to the crucifix (a cross that includes a corpus, usually a three-dimensional ...

Roof



A roof is the top covering of a building, including all materials and constructions necessary to support it on the walls of the building or on uprights, providing protection against rain, snow, sunlight, extremes of temperature, and wind. A ...

The Role of Context for Evaluation

Building Material



Accessibility Description 1

Building material is material used for construction. Many naturally occurring substances, such as clay, rocks, sand and wood, even twigs and leaves, have been used to construct buildings. Apart from naturally occurring materials, ...

Christian Cross



Accessibility Description 2

The Christian cross, seen as a representation of the instrument of the crucifixion of Jesus, is the best-known symbol of Christianity. It is related to the crucifix (a cross that includes a corpus, usually a three-dimensional ...

Roof



Accessibility Description 3

A roof is the top covering of a building, including all materials and constructions necessary to support it on the walls of the building or on uprights, providing protection against rain, snow, sunlight, extremes of temperature, and wind. A ...

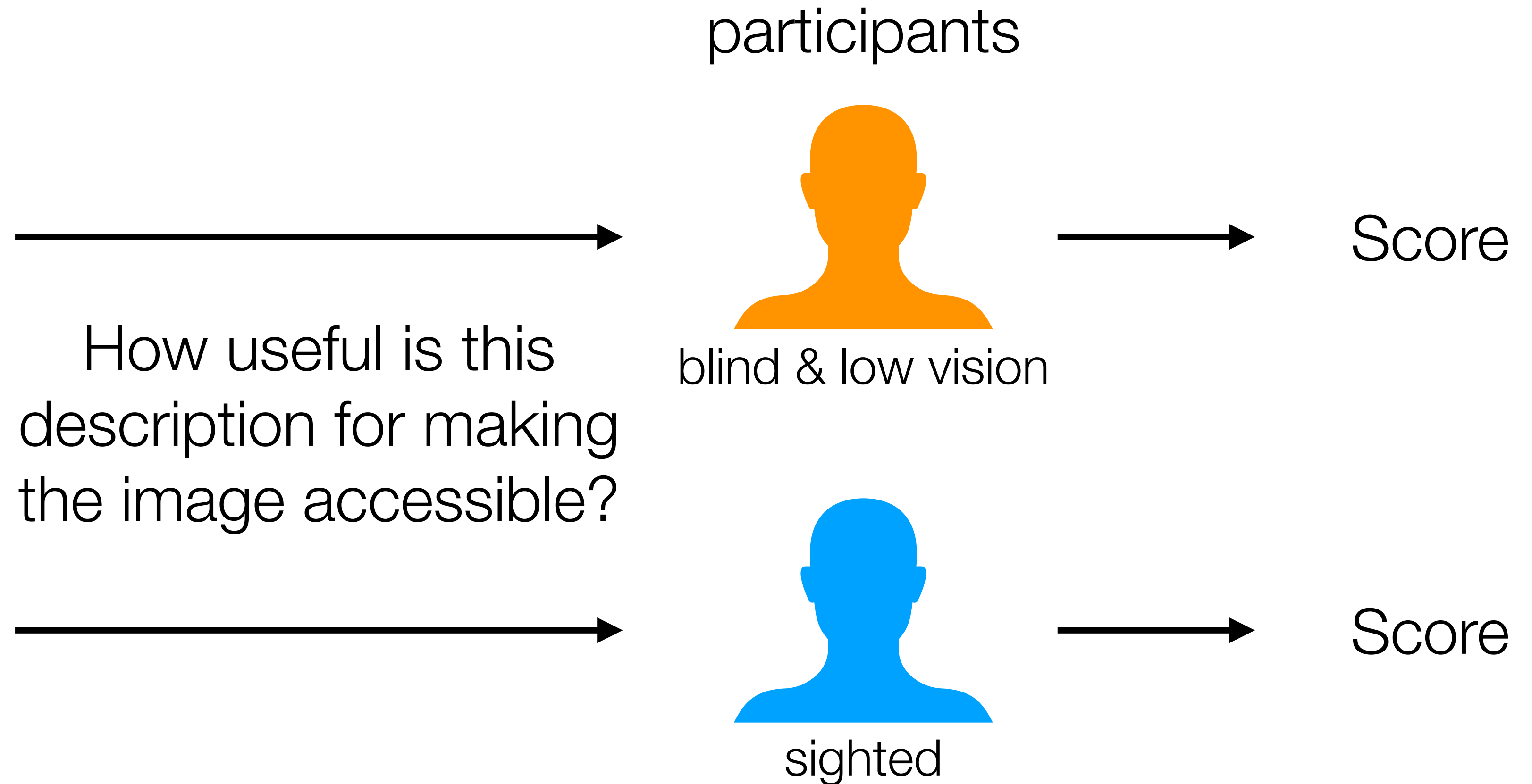
The Role of Context for Evaluation

Building Material

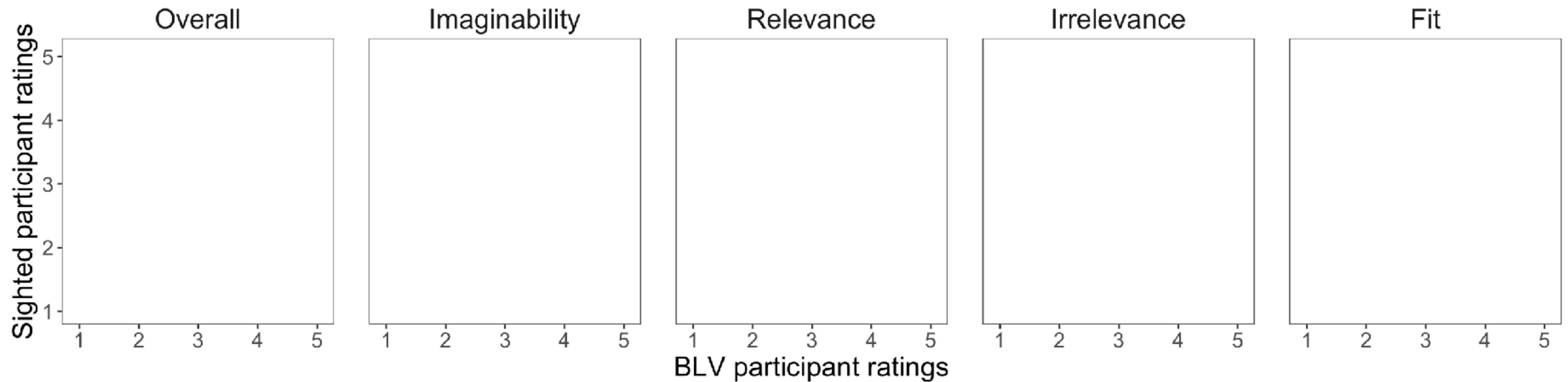


Accessibility Description 1

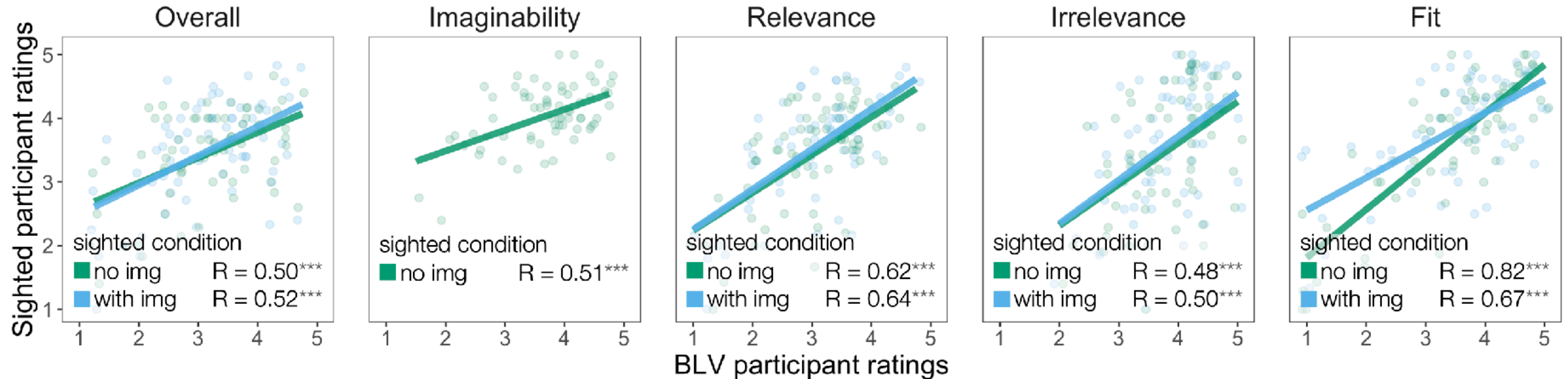
Building material is material used for construction. Many naturally occurring substances, such as clay, rocks, sand and wood, even twigs and leaves, have been used to construct buildings. Apart from naturally occurring materials, ...



The Role of Context for Evaluation

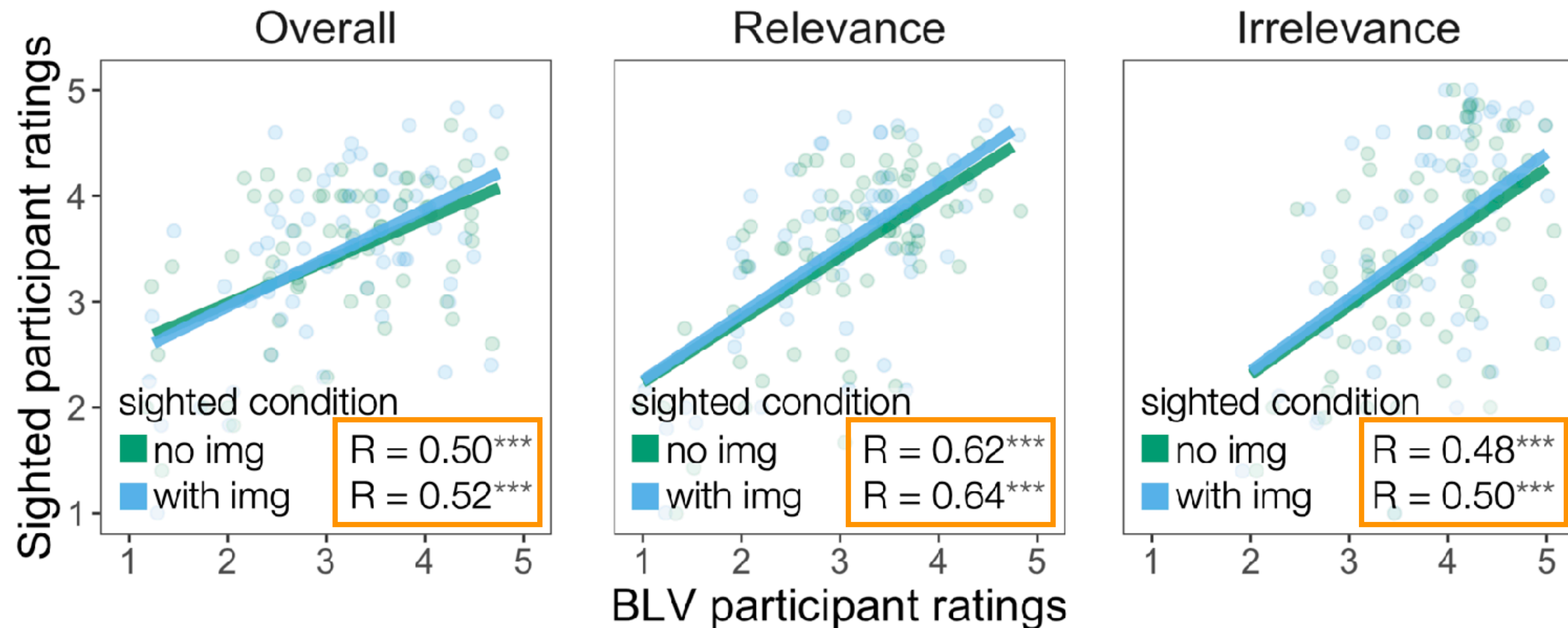


The Role of Context for Evaluation



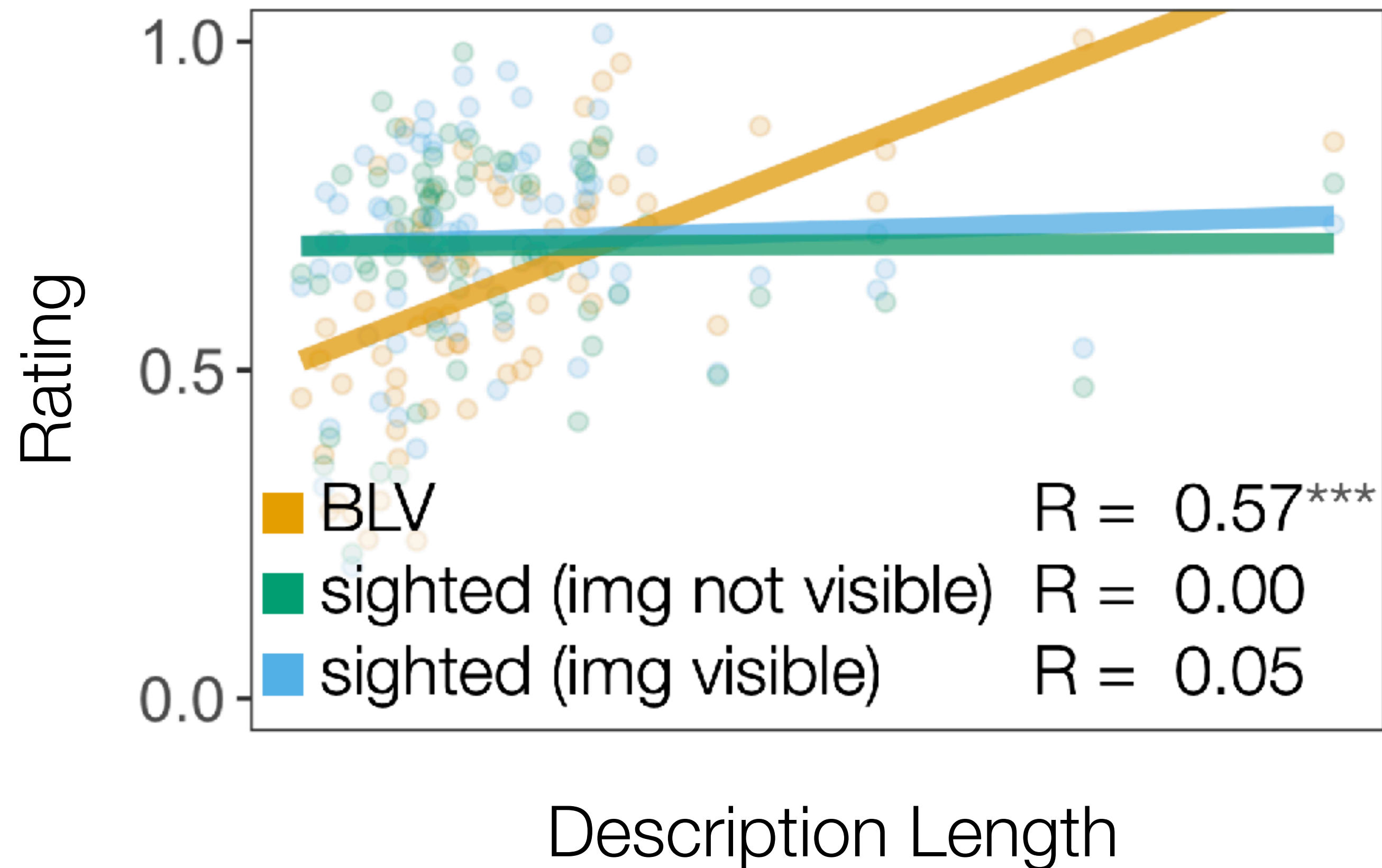
The Role of Context for Evaluation

Hiding the image from sighted participants doesn't improve their alignment with BLV participant ratings.

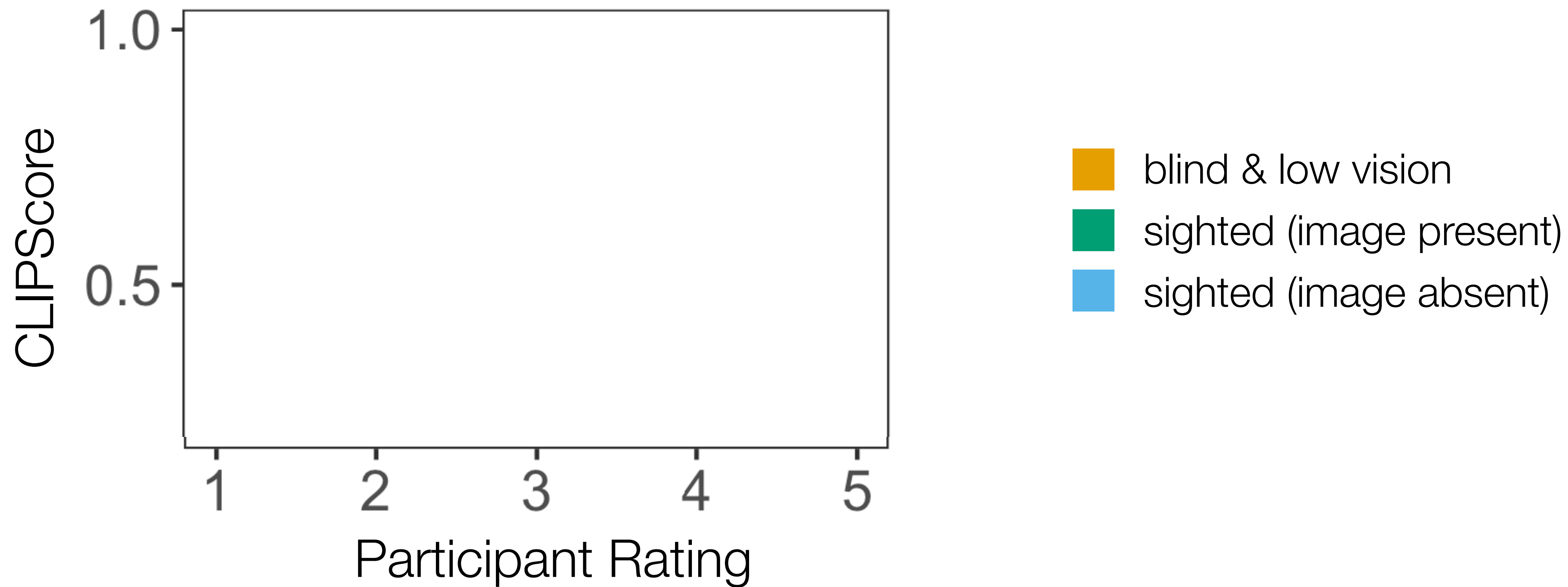


The Role of Context for Evaluation

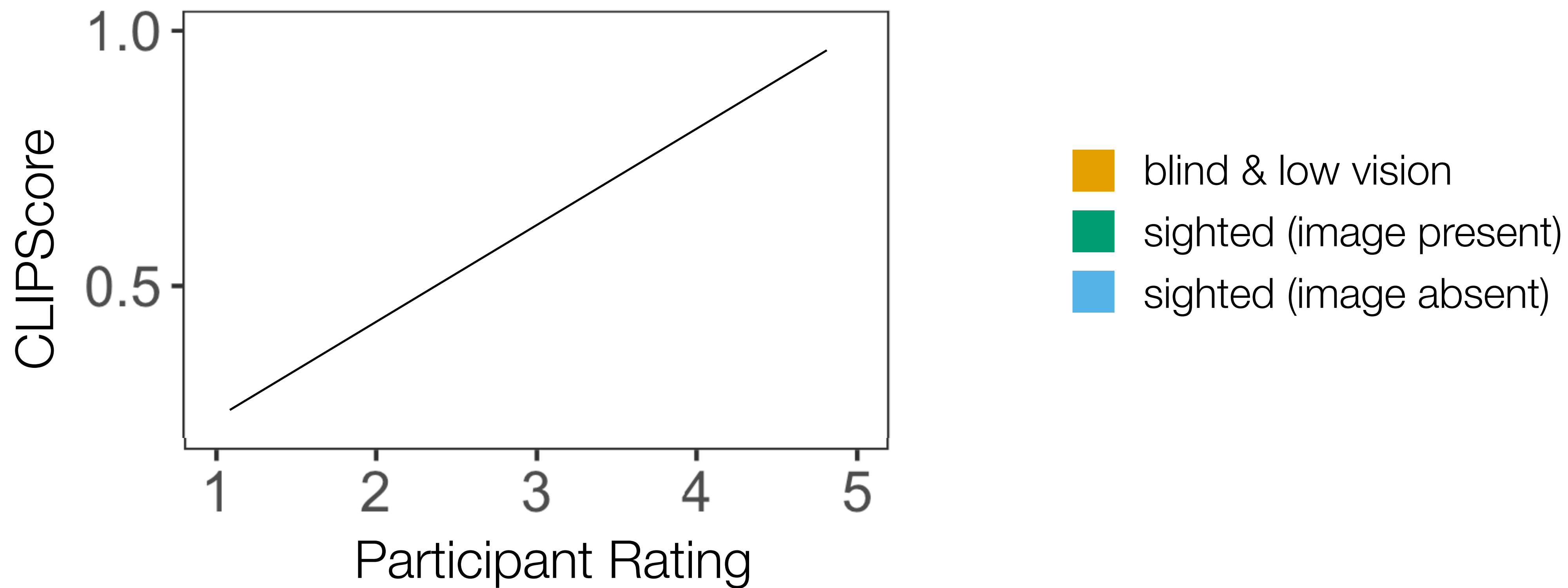
BLV participant ratings (but **not** sighted participant ratings) show a strong correlation with description length



The Role of Context for Evaluation

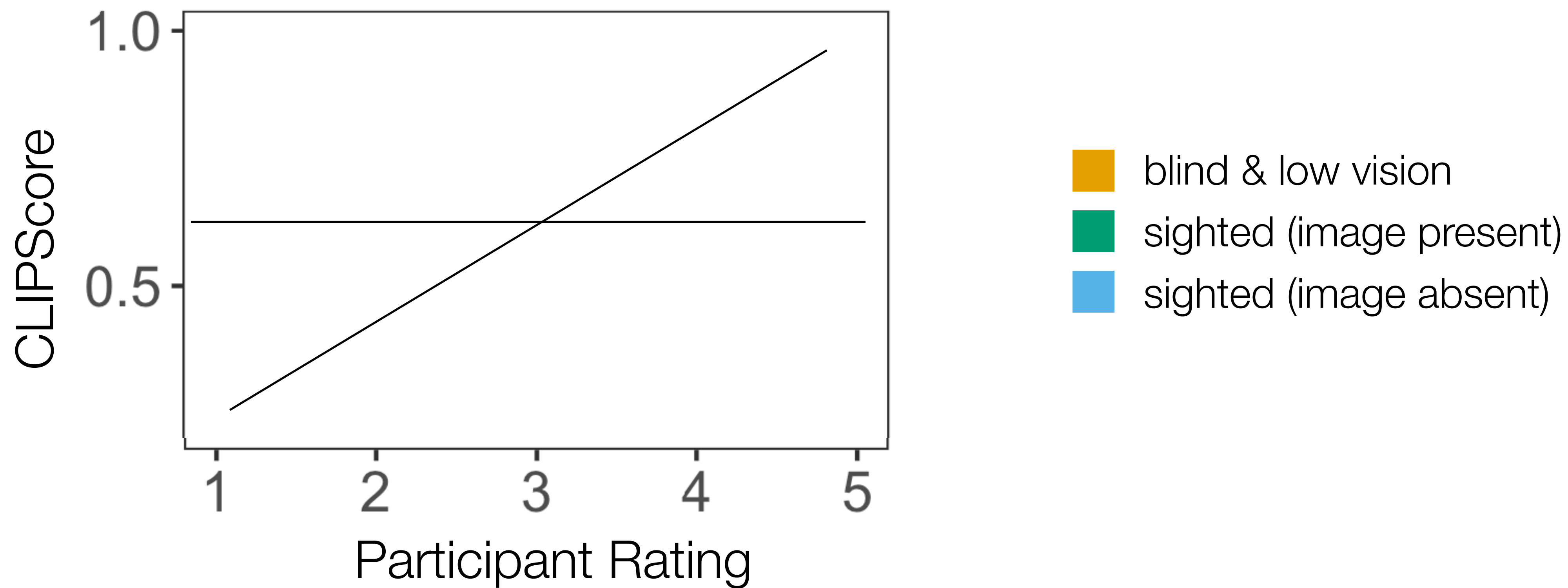


The Role of Context for Evaluation



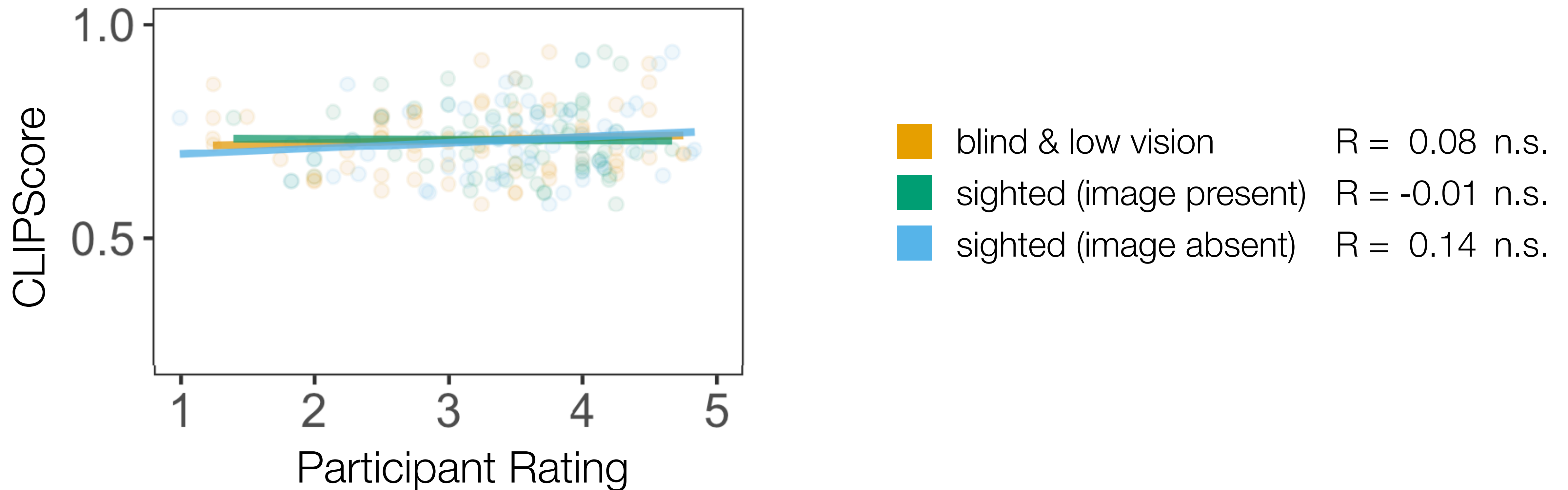
Context matters for image descriptions for accessibility: Challenges for referenceless evaluation metrics

The Role of Context for Evaluation



The Role of Context for Evaluation

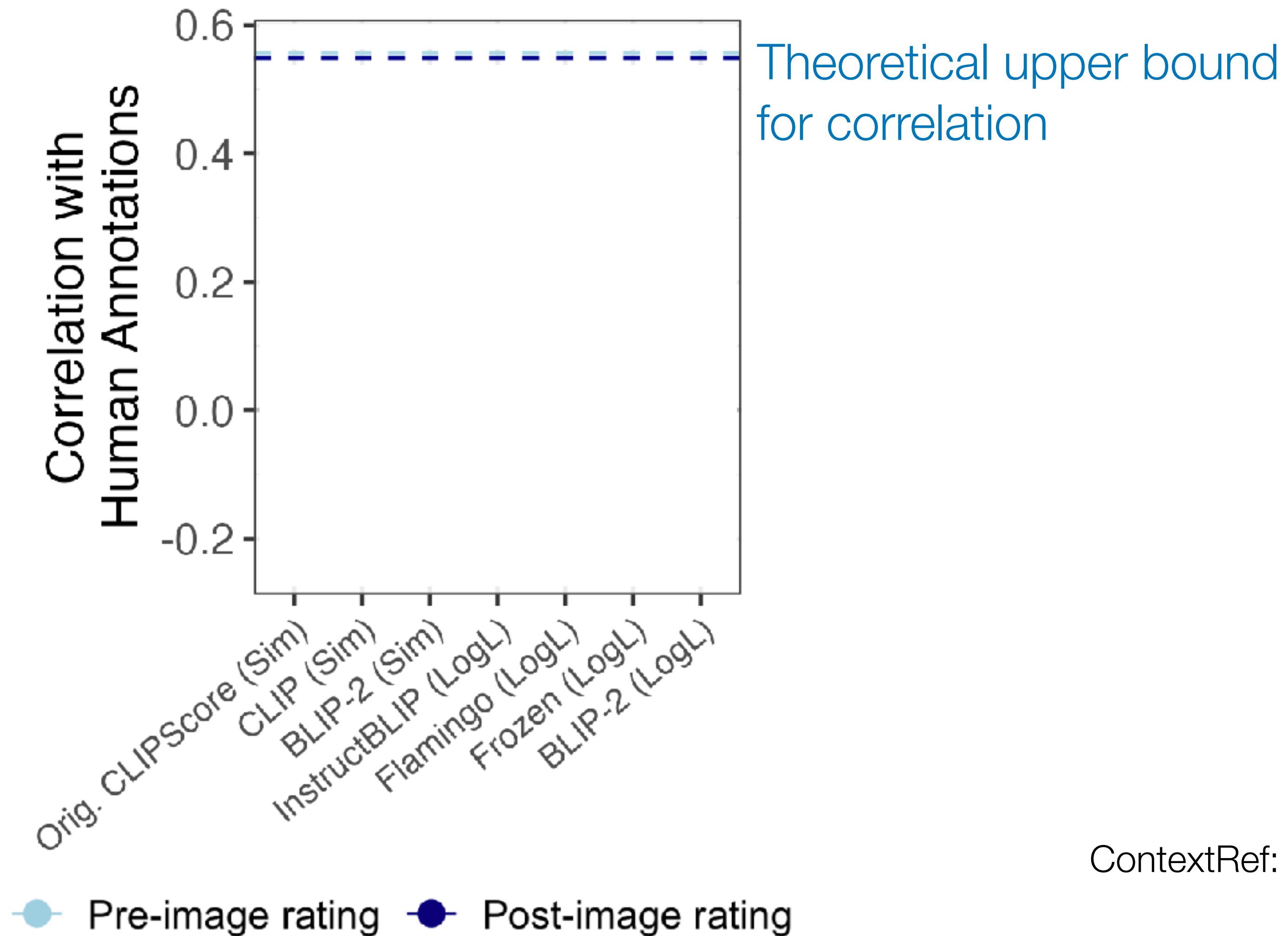
No correlation between CLIPScore and human ratings



Models for Human-like Image Description Evaluation

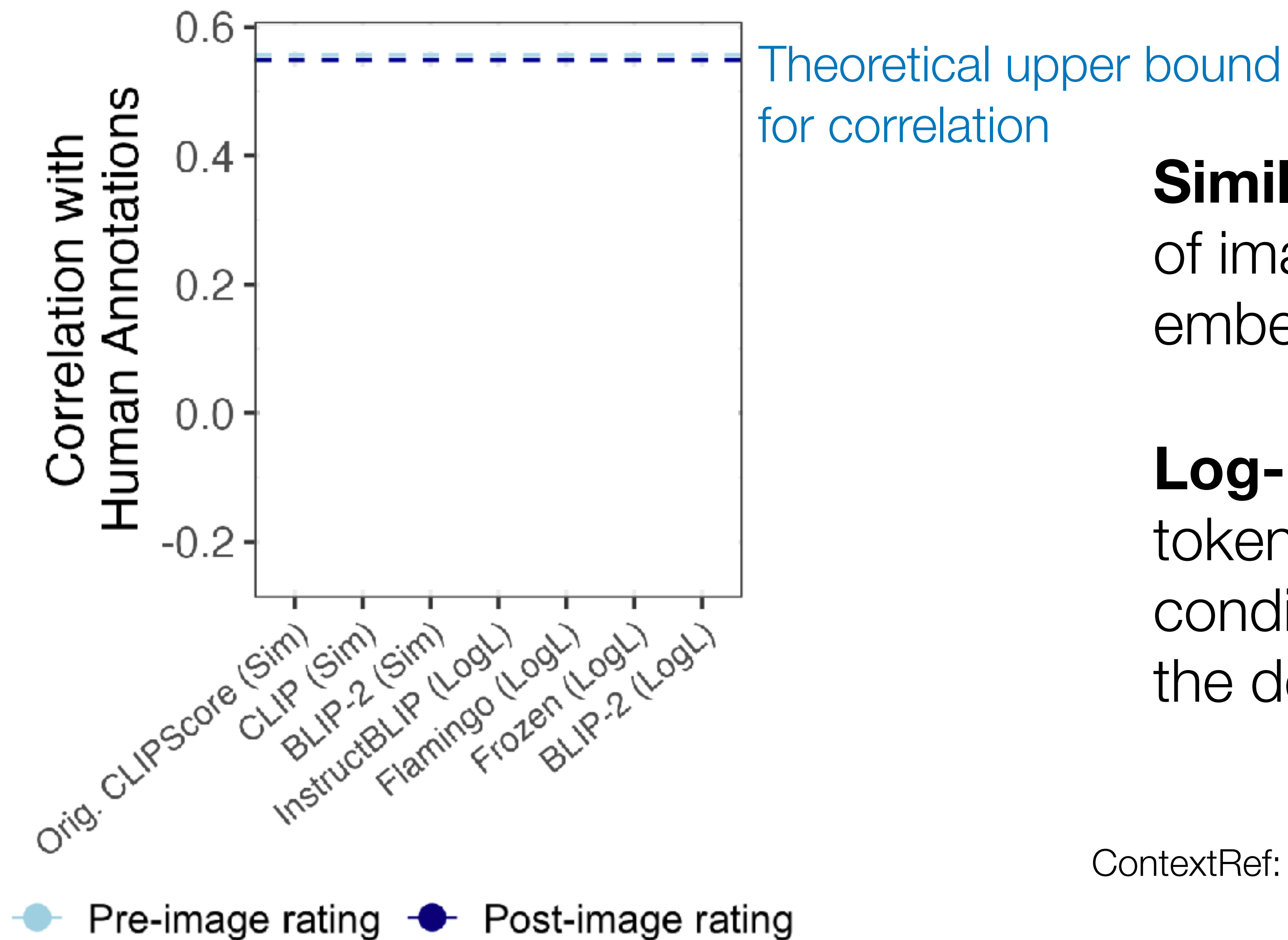
Models for Human-like Image Description Evaluation

Promising correlations with human raters



Models for Human-like Image Description Evaluation

Promising correlations with human raters

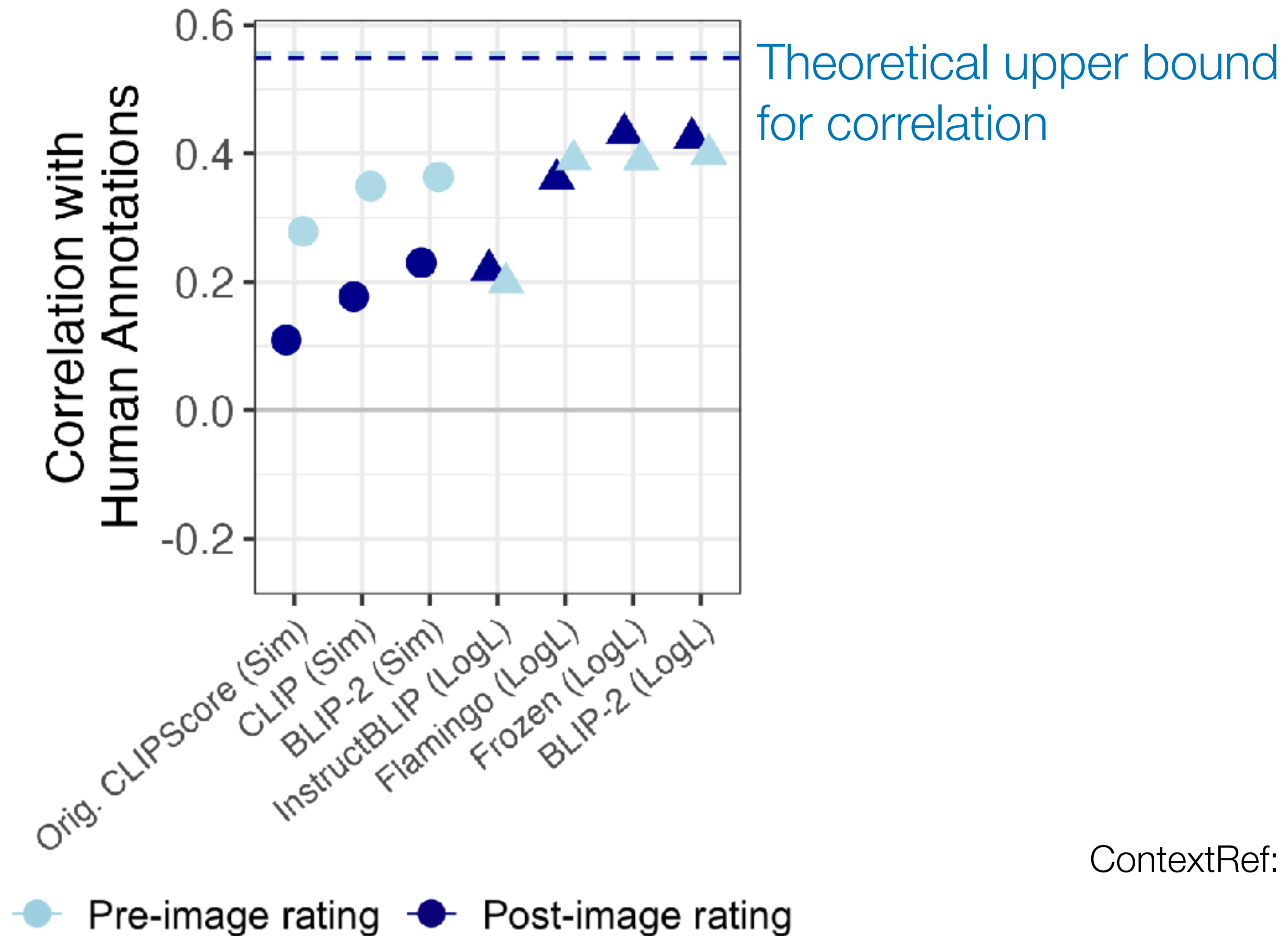


Similarity-based metrics: Cosine similarity of image, description, and context in embedding space

Log-likelihood metrics: LM's average per-token log-likelihood of the description, conditioned on the image, the context, and the description being high-quality

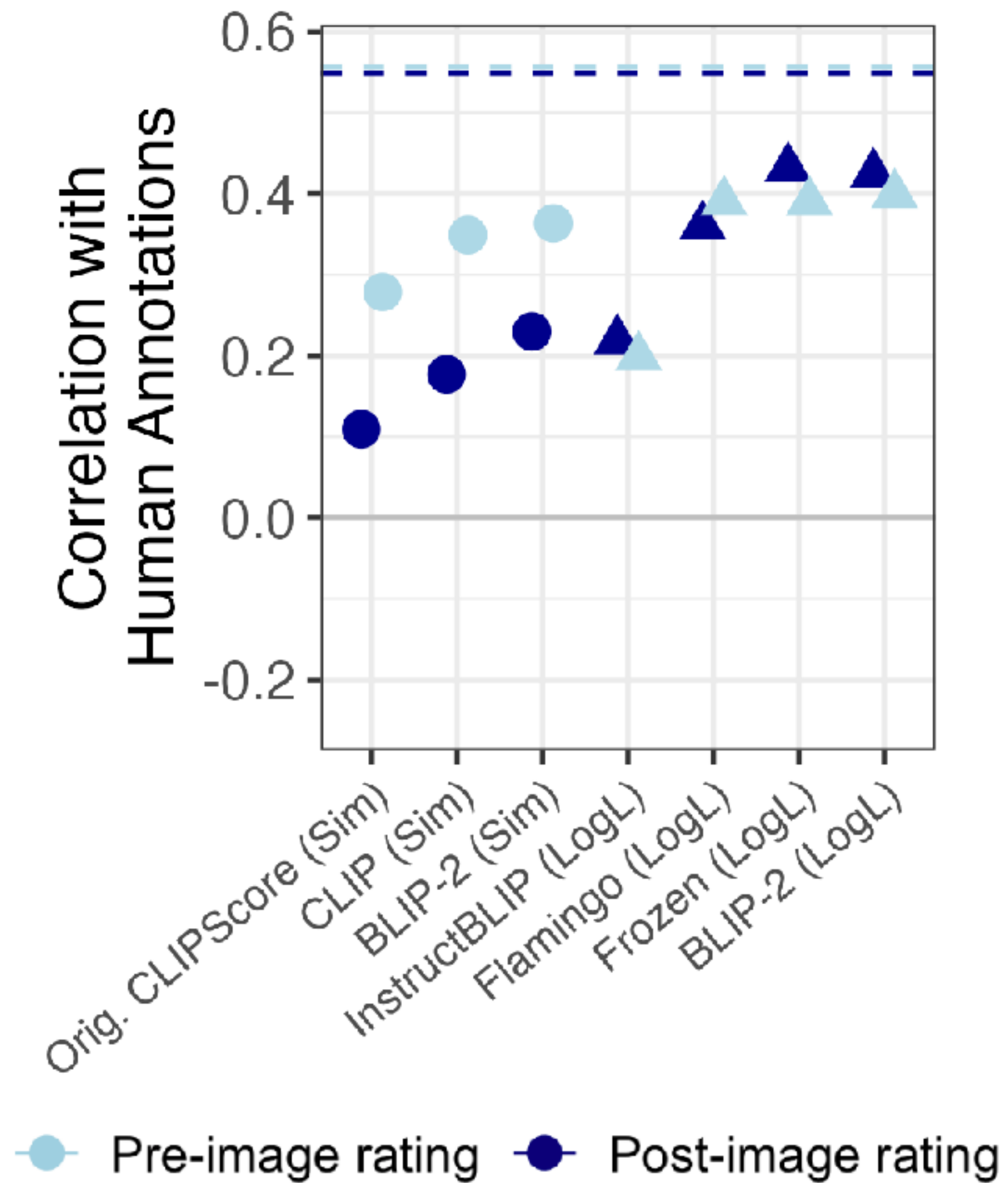
Models for Human-like Image Description Evaluation

Promising correlations with
human raters



Models for Human-like Image Description Evaluation

Promising correlations with human raters



Data augmentations uncover unreliable behavior in all models

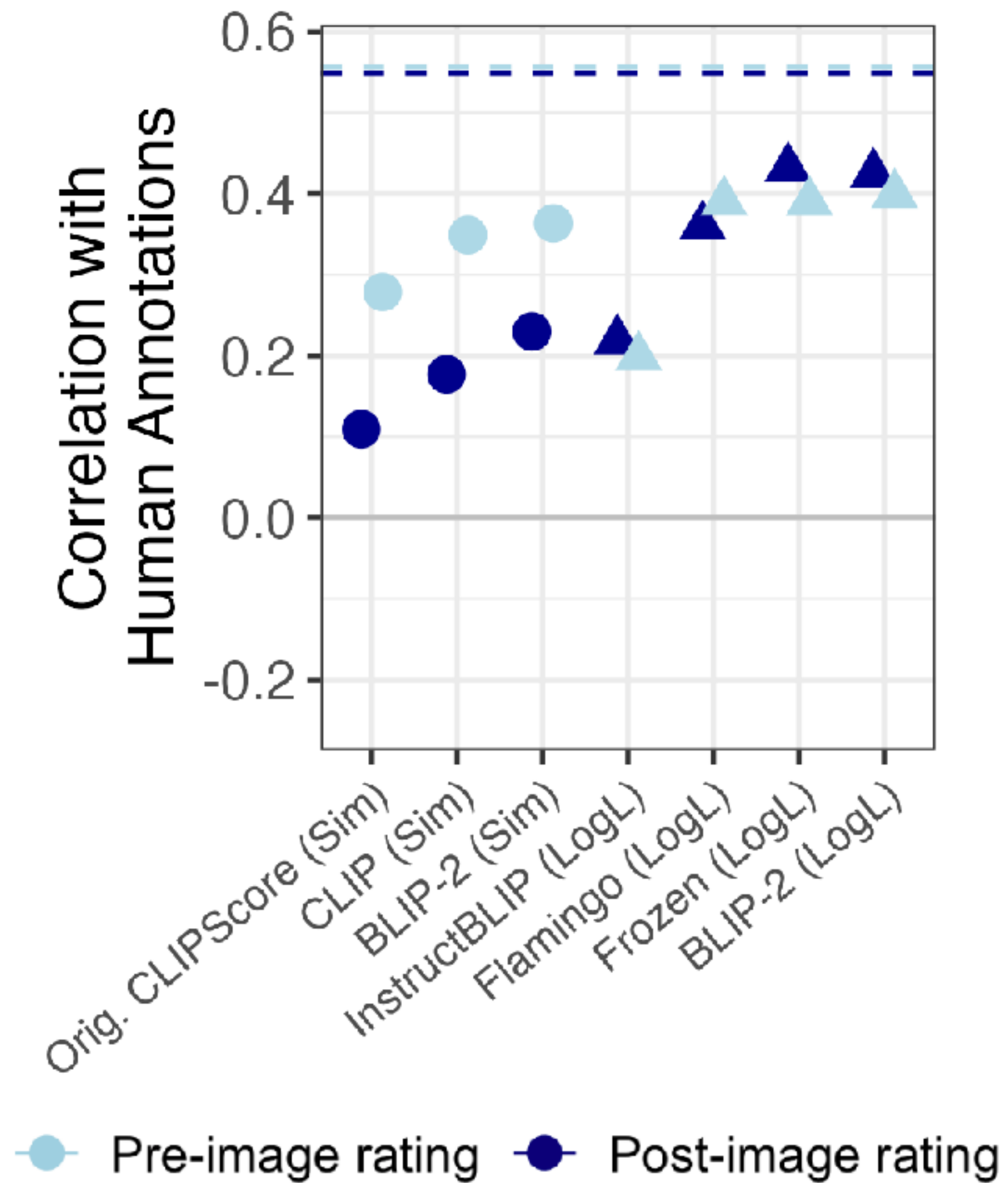


Original: a dog with big ears

Good

Models for Human-like Image Description Evaluation

Promising correlations with human raters



Data augmentations uncover unreliable behavior in all models



Original: a dog with big ears

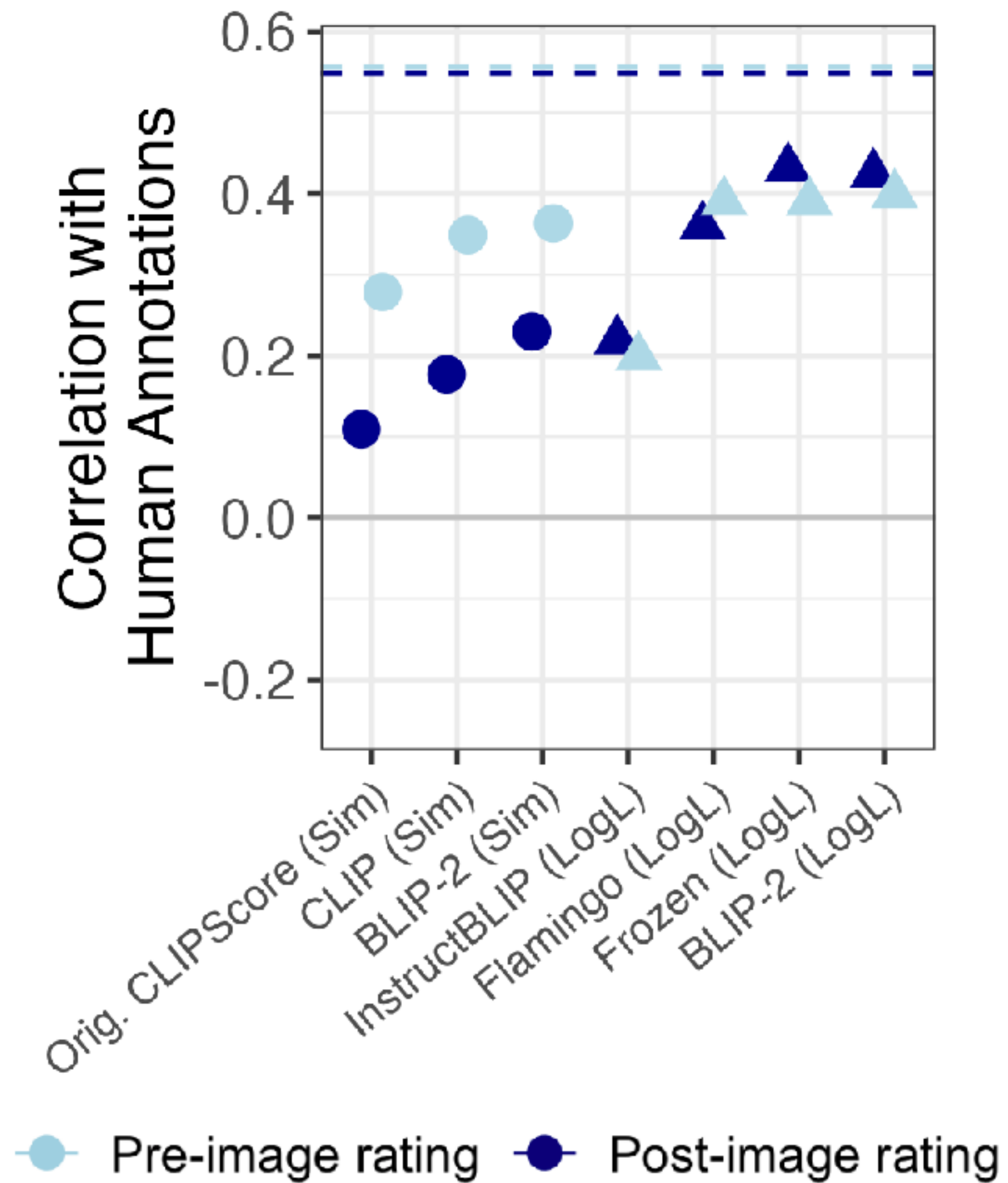
Good

Exact repetition: a dog with big ears
ears a dog with big ears

Excellent!

Models for Human-like Image Description Evaluation

Promising correlations with human raters



Data augmentations uncover unreliable behavior in all models



Original: a dog with big ears

Good

Exact repetition: a dog with big ears
a dog with big ears

Excellent!

Irrelevance: a dog with big ears
Elephants are the largest existing land animals

Excellent!

Models for Human-like Image Description Evaluation

Data augmentations uncover
unreliable behavior in all models

Even state-of-the-art
models are still far
from approximating
real human data.



Original: a dog with big ears

Good

Exact repetition: a dog with big
ears a dog with big ears

Excellent!

Irrelevance: a dog with big
ears Elephants are the largest
existing land animals

Excellent!

Image-based text generation depends on ...

1 the **image-based text's** communicative goal.

→ description ≠ caption

A sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. Part of the root system are the taproots and lateral roots. The taproot refers to the the central root and the lateral roots are the smaller side roots that ...

A diagram of the anatomy of a plant with labels of structural parts of the plant and the roots.

2 the **image's** communicative goal.

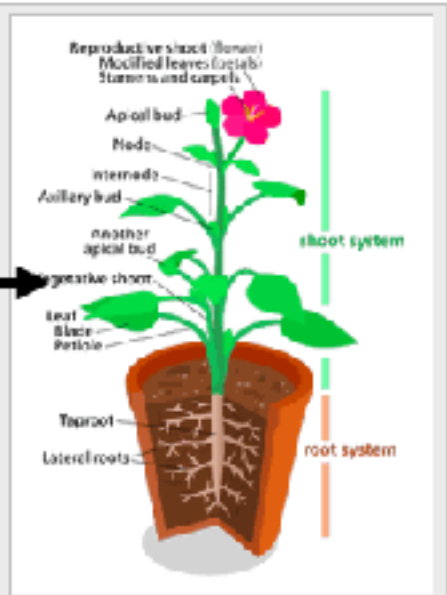
→ article context determines information selection



Multimodal Pedagogy

Multimodal pedagogy is an approach to the teaching of writing that implements different modes of communication.^{[1][2]} Multimodality refers to the use of visual, aural, linguistic, spatial, and gestural modes in differing pieces of media, each necessary to convey meaning.^[3] The visual mode conveys meaning via the use of visual elements such as images, diagrams, and color. The aural mode refers to sound in the form of audio recordings or spoken language. The spatial mode refers to physical movement or layout in space. The gestural mode refers to physical movement or layout in space. Multimodal text is characterized by the combination of any two or more modes to convey meaning.^[4] Multimodality as a term was coined in the late 20th century,^[5] but its use predates its naming, with it being used as early as Egyptian hieroglyphs and classical rhetoric.^[7] Compositionists and writing theorists have been exploring how the five modes of communication interact with each other and how multimodality can be used in the teaching of writing since the 20th century.^[6]

An educational sketch of a plant, illustrating the shoot system which is above the soil, and the root system which is below the soil. The different parts of the plant are labeled to point out, e.g., the plant's reproductive shoot (i.e., the flower) or the lateral roots.



Example of complementing linguistic and visual information leading to learning benefits over unimodal approaches.

What **can** we say
about an image?



What **should** we say
about an image?

WIKIPEDIA

The Free Encyclopedia

Plant Anatomy

Plant anatomy or **Phytotomy** is the general term for the study of the internal **structure** of **plants**. Originally it included **plant morphology**, the description of the physical form and external structure of plants, but since the mid-20th century plant anatomy has been considered a separate field referring only to internal plant structure.^{[1][2]} Plant anatomy is now frequently investigated at the **cellular level**, and often involves the sectioning of **tissues** and **microscopy**.^[3]

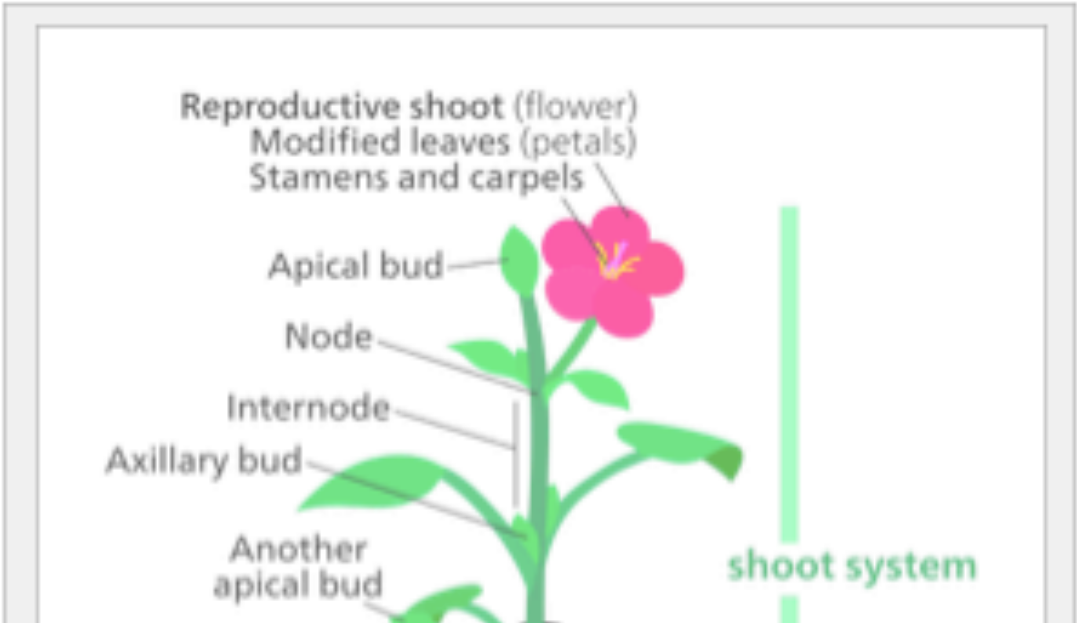


Image-based text generation depends on ...



the image-based **text**'s communicative goal.

→ description \neq caption

the **image**'s communicative goal.

→ article context determines information selection

What **can** we say
about an image?



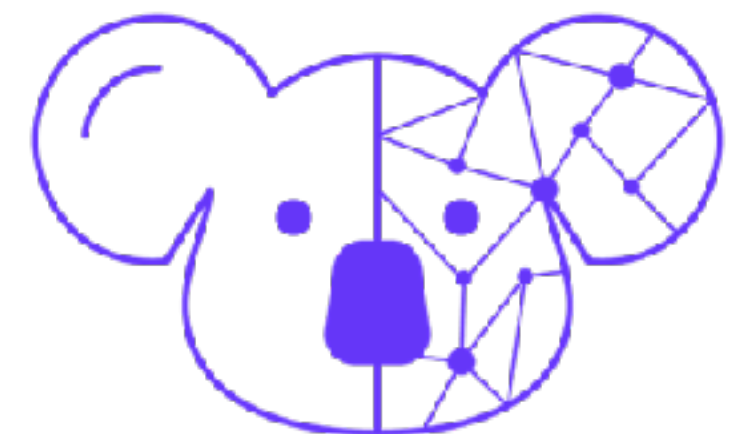
What **should** we say
about an image?

Frontiers in Generative AI for Nonvisual Accessibility

- Development of **local** (privacy-preserving, internet-independent) systems likely requires the development of smaller specialized systems which require more task-specific supervision.
- **Well-calibrated** systems (including hallucinations) for which we can anticipate error behavior and achieve pragmatic inference alignment (relevance, informativity, length) are especially consequential in situations where verification of the ground-truth is hard.

ekreiss@ucla.edu

@elisakreiss



Coalas Lab