



LECTURE 1

GEOSPATIAL REPRESENTATION LEARNING

Introduction to Geospatial Representations

PRESS – OR SPACE. →

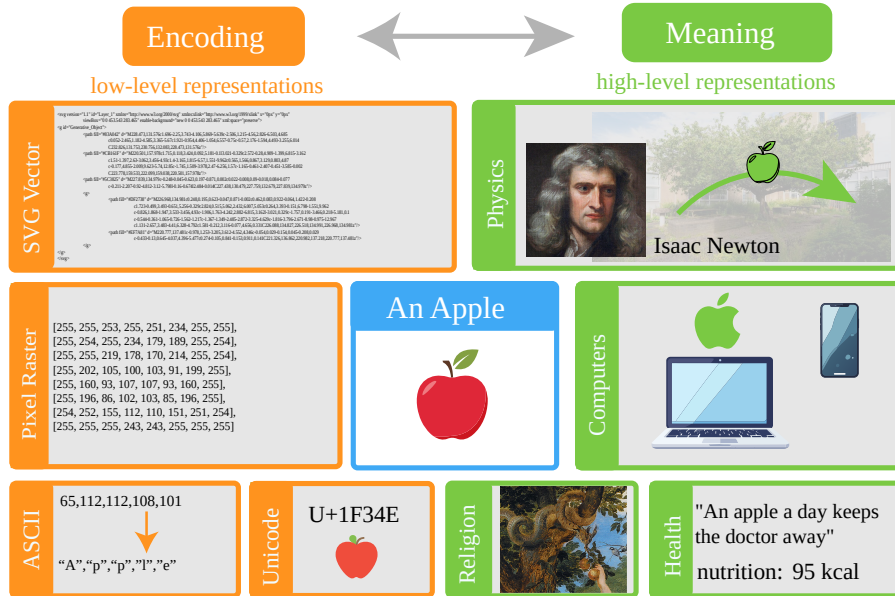
Representations of our World



[Apple](#) by Fir0002/Flagstaffotos license CC-BY-NC, Earth added via ChatGPT

Describe an Apple.

Meaning and Encoding of Representations



The same object can be encoded in many ways. A good encoding representation makes useful structure visible and computation easier.

But representation is not only encoding. Representation is also meaning.

Choosing Right Encoding Representations

Examples

Images

255	255	255	255	255	55	55	255	255	255	255	255
255	255	255	255	55	55	255	55	55	255	255	255
255	255	255	70	170	170	70	170	170	70	255	255
255	255	70	200	200	200	200	200	200	200	70	255
255	70	200	200	200	200	200	200	200	200	200	70
255	70	200	200	200	200	200	200	200	200	200	70
255	70	200	200	200	200	200	200	200	200	200	70
255	70	200	200	200	200	200	200	200	200	200	70
255	255	70	200	200	200	200	200	200	200	70	255
255	255	70	200	200	200	200	200	200	200	70	255
255	255	255	70	200	200	200	200	200	70	255	255
255	255	255	255	70	70	70	70	70	70	255	255

Math

XLVII + LXXVIII = ?

$$47 + 78 = 125$$

Places

COORDINATES

50°43'57.73" N,
7°06'16.63" E

PLACE

**Hofgarten, Bonn,
Germany**

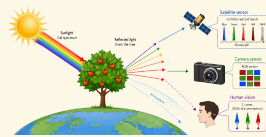
Lecture Outline

Representations



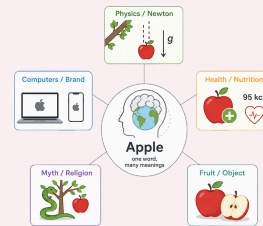
What are representations?

Vision



How vision encodes the world

Language



How text becomes meaning

Spatial

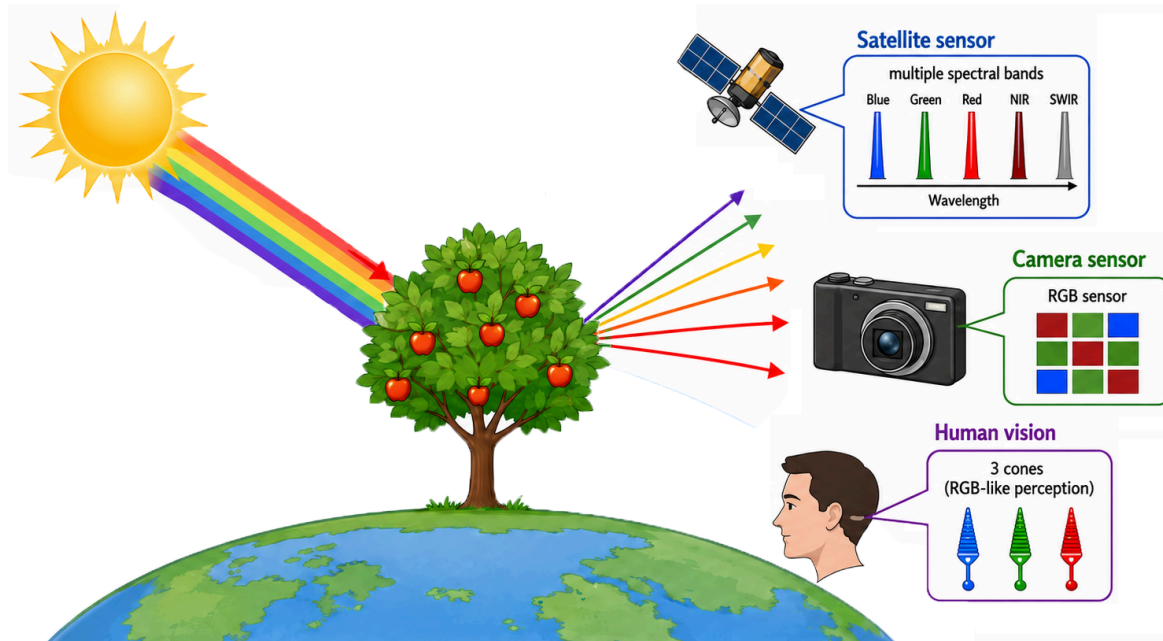


How places carry spatial meaning

Visual Representations

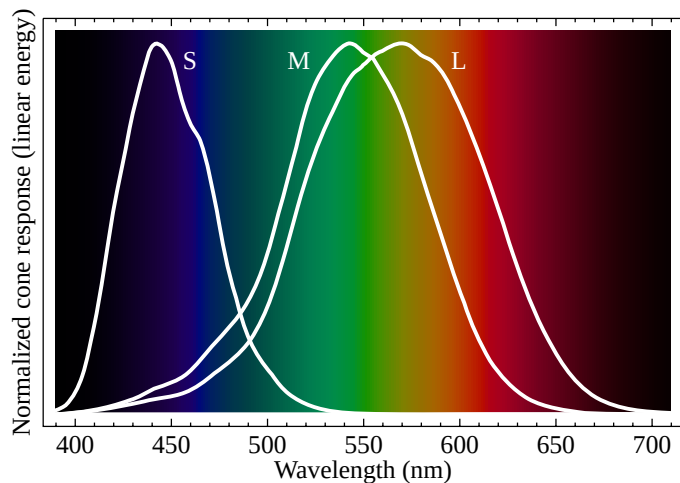
We perceive our world in representations

Example: Vision

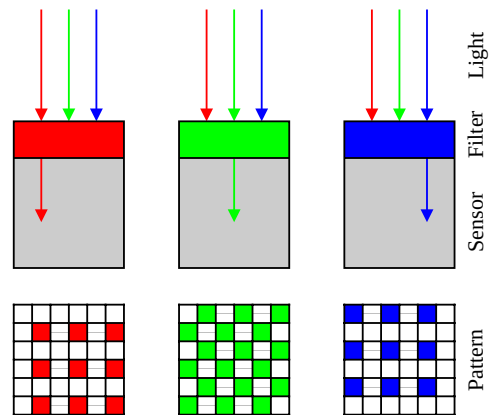


Engineering suitable Representations.

Sensitivity of cone cells (S, M, L) in human eyes

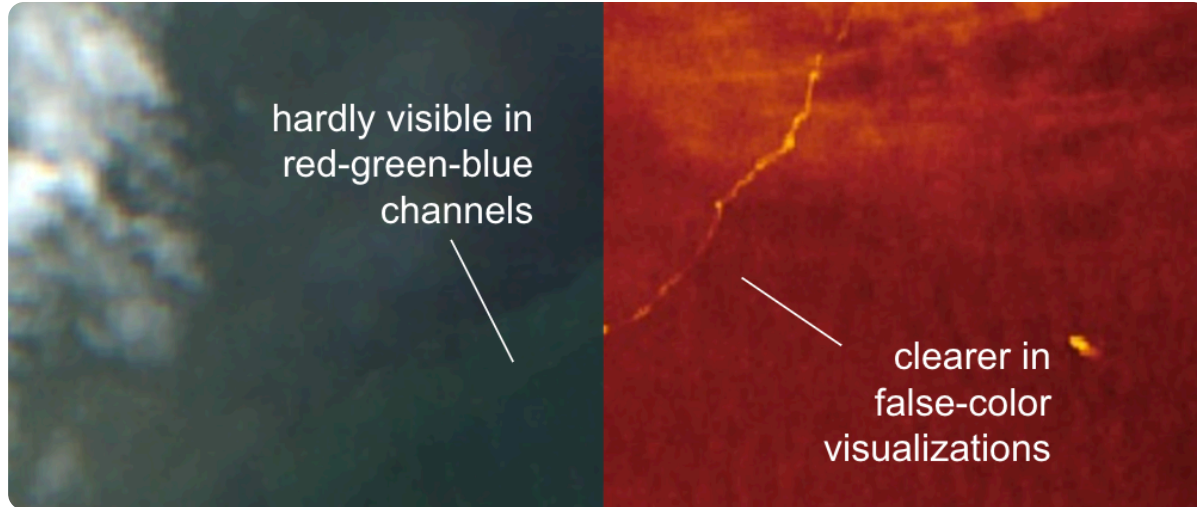


Bayer Filter on Cameras



Cameras are built to mimic our eyes: More green pixels in CCD camera mimics the color sensitivity in our eyes.

Spectral information can make patterns visible.



Example: plastic litter and marine debris (Sentinel-2)

The representation depends on the Problem

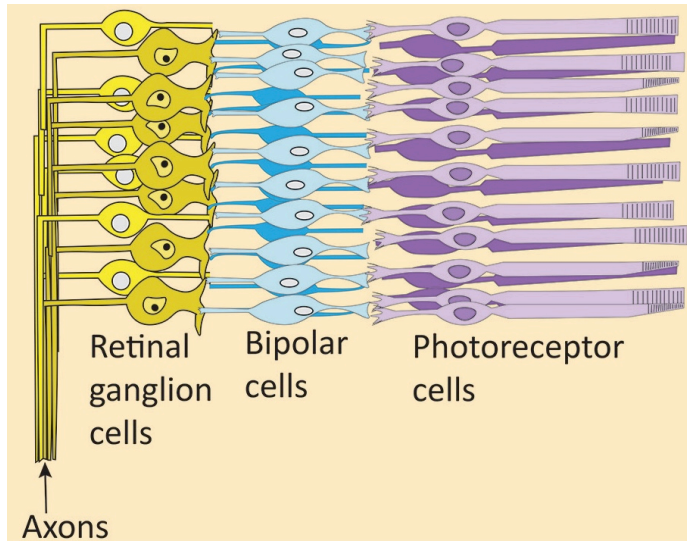


A useful representation depends on the problem we are trying to solve.

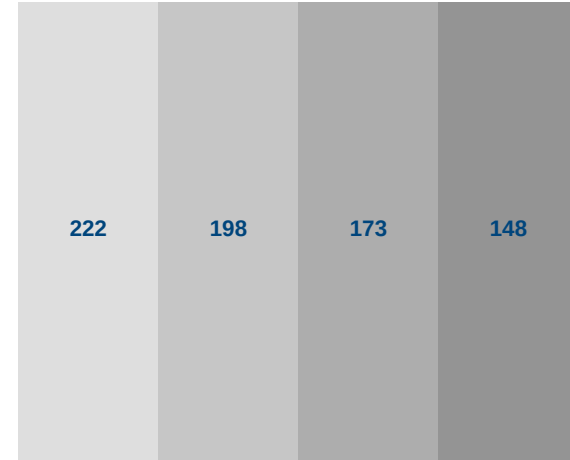
Bees see ultraviolet patterns that humans miss, because their visual system is tuned to different tasks.

Our eyes have edge detectors

Retinal cells implement an edge filter through lateral inhibition.



Mach Bands - Optical Illusion (Mach, 1865)



Our eyes are not objective, they capture an edge-enhanced representation of the world.

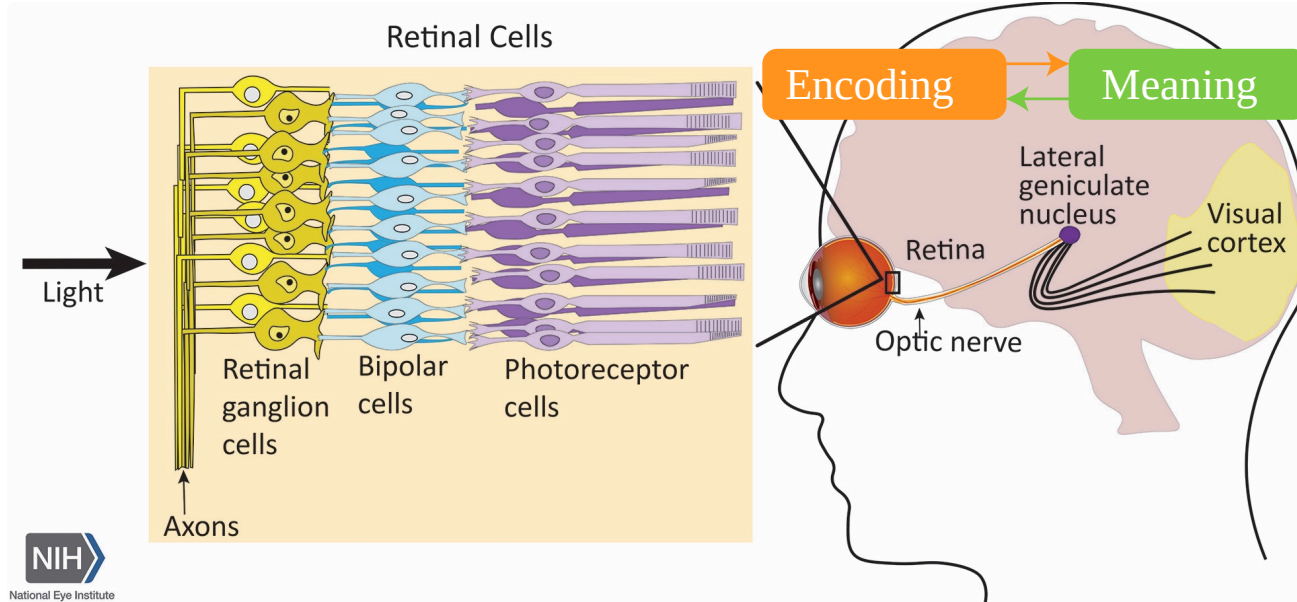
Shapley, R. M., & Tolhurst, D. J. (1973). Edge detectors in human vision. *Journal of Physiology*, 229, 165-183.

Selective Vision: The Monkey Business Illusion



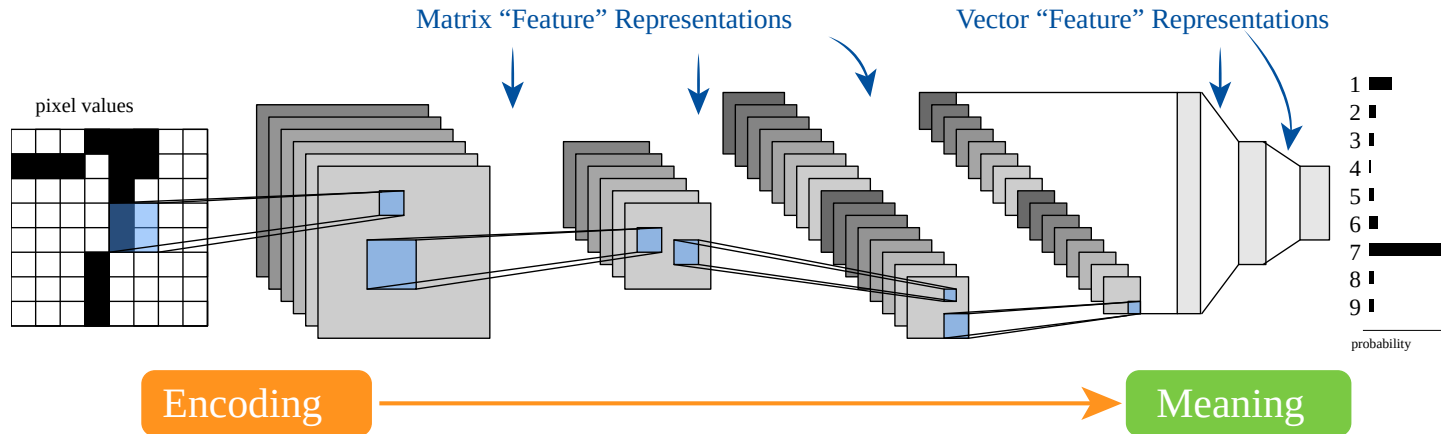
Simons, D. J. (2010). Monkeying around with the gorillas in our midst: Familiarity with an inattentional-blindness task does not improve the detection of unexpected events. *i-Perception*. <https://doi.org/10.1068/i0386>

Our mind selects what we perceive



From Encoding to Meaning

Feature Learning in Neural Networks



Excuse: Dimensionality Reduction

We often need to visualize high-dimensional data on a lower-dimensional space, as here with 3D data on a 2D screen.

Principal Component Analysis (PCA) linearly rotates high-dimensional data points so that a 2D projection best captures the variance (i.e., spread) of the data

Rotate to PCA projection

Reset

Show projections

t-distributed Stochastic Neighborhood Embedding (t-SNE) moves data points onto a 2D plane and tries to preserve the local neighborhood relation between points.

Run t-SNE projection

Reset

Show trajectories

Feature Representations in LeNet (2000)

Supervised Learning: From Labels to Loss

Learning from labeled images

- **Input:** a 28×28 image of a handwritten digit
- **Known truth:** the human-provided label 7
- **Prediction:** probabilities over the ten digit classes
- **Learning:** iteratively adjust the model to minimize the prediction error
- **Goal:** make the probability for the true label 7 as high as possible

```
image pixels → CNN representation → class  
probabilities → label
```

How do we measure whether the prediction is correct or incorrect?

Classification example: Cross-Entropy Loss

$$\mathcal{L}_{CE}(y, \hat{y}) = - \sum_{c=1}^C y_c \log(\hat{y}_c)$$

For a one-hot target label:

$$\mathcal{L}_{CE} = - \log(\hat{y}_{\text{correct}})$$

Correct class: 7

- If $p(7) = 0.90$ → low loss
- If $p(7) = 0.10$ → high loss

Cross entropy rewards high probability on the correct class and penalizes confidence in wrong classes.

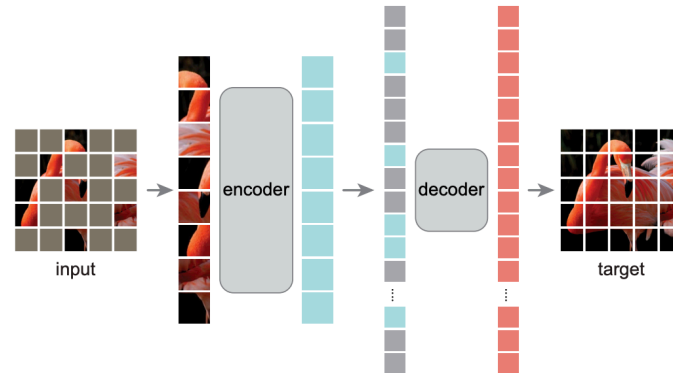
Self-supervised Learning

Contrastive Learning

Identify the corresponding object among two different views of the same object.

Masked Autoencoding

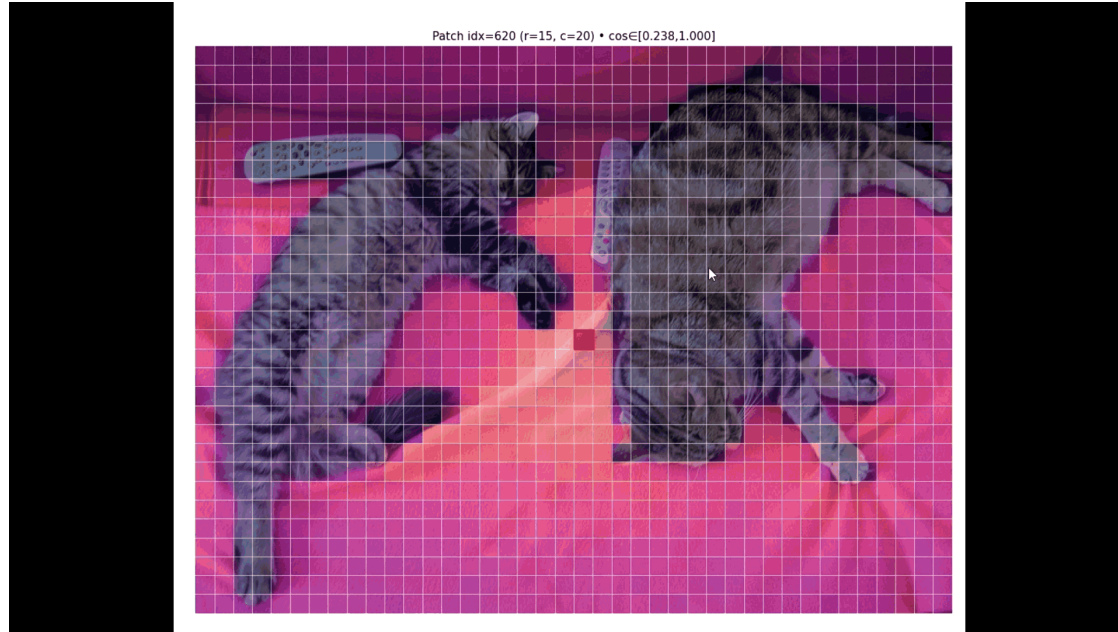
Train the model to correctly inpaint a masked portion of an image.



Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., & Joulin, A. (2020). [Unsupervised Learning of Visual Features by Contrasting Cluster Assignments](https://arxiv.org/abs/2006.05759). NeurIPS. Code: [facebookresearch/swav](https://github.com/facebookresearch/swav).

He, K., Chen, X., Xie, S., Li, Y., Dollár, P., & Girshick, R. (2022). [Masked autoencoders are scalable vision learners](https://arxiv.org/abs/2203.09788). CVPR, 16000-16009.

High-level Visual Features (Dinov3, 2025)



Demo: [devMuniz02/DINOV3-Interactive-Patch-Cosine-Similarity](#)

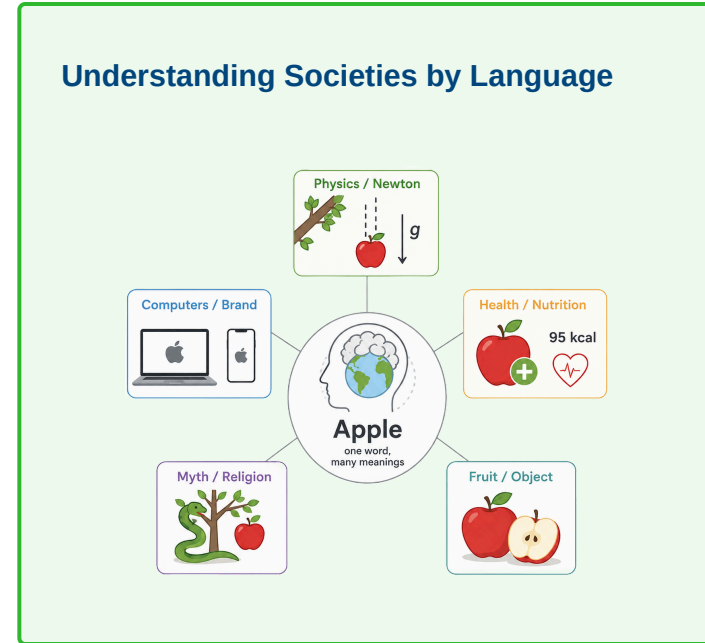
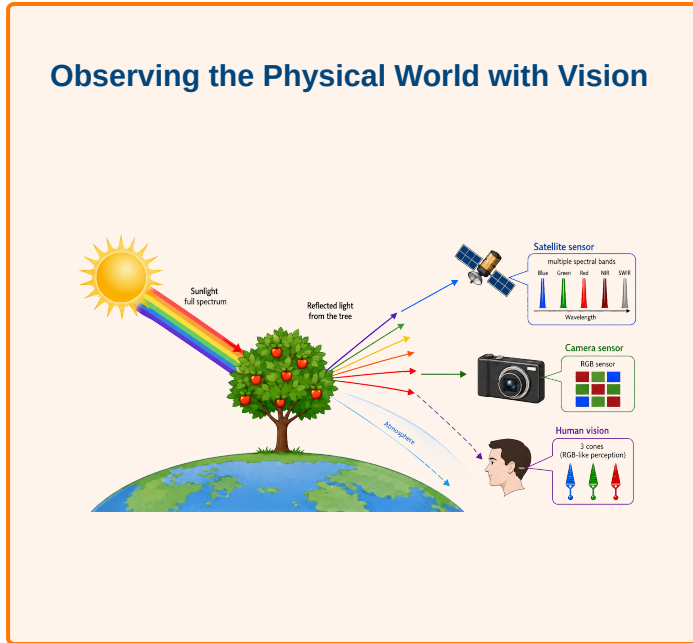
DINOV3 model: [Siméoni, O., Vo, H. V., Seitzer, M., Baldassarre, F., Oquab, M., Jose, C., Khalidov, V., Szafraniec, M., Yi, S., Ramamonjisoa, M., et al. \(2025\). DINOV3. arXiv:2508.10104.](#)

High-level Visual Features (SAM 2, 2024)

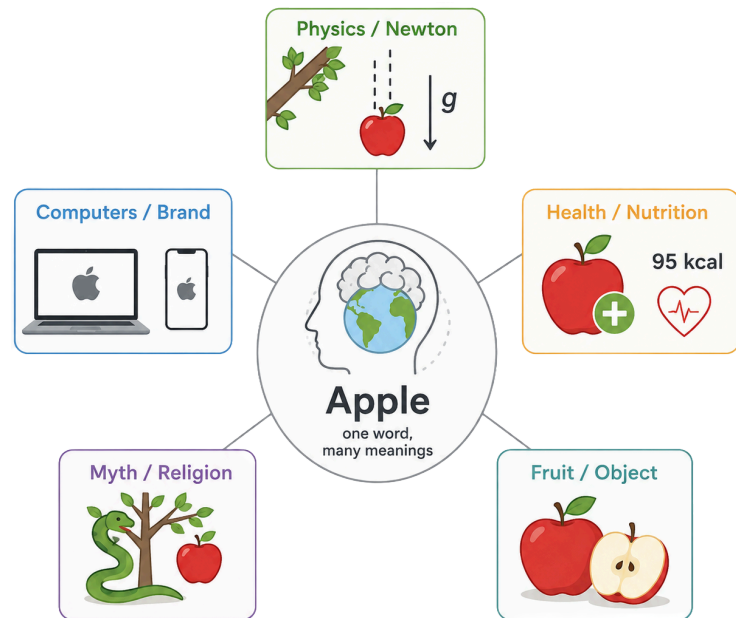


Try the Segment Anything Model (SAM) 2 demo

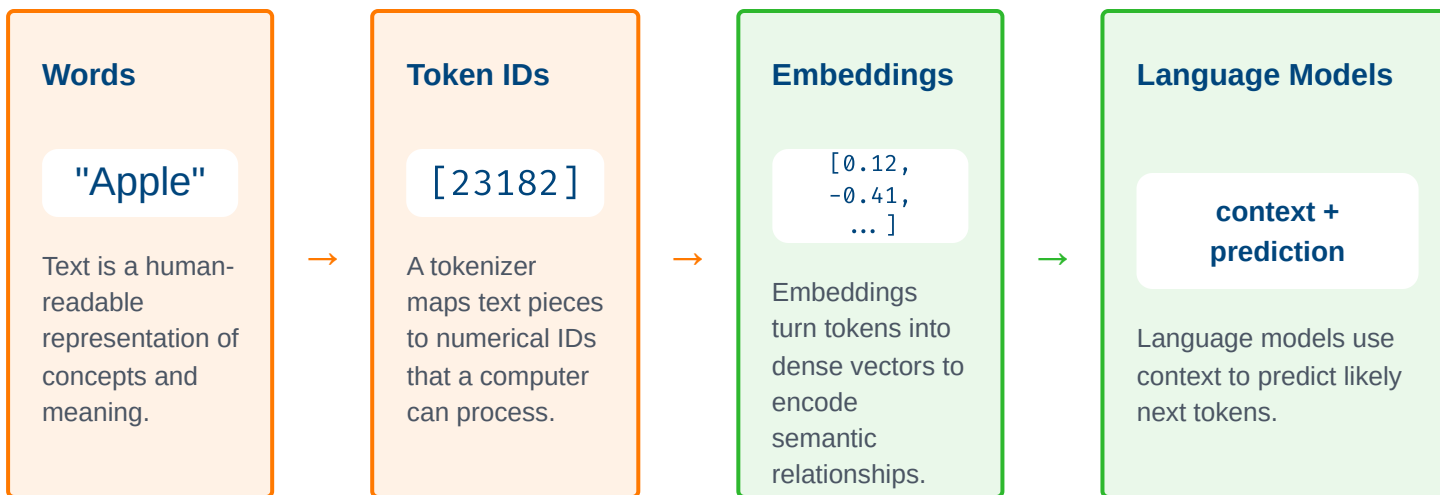
From Physical to Societal Representations



Language Representations



From Text Encoding to Language Models



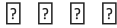





Text and language captures the human world: machines move from low-level encodings of words and tokens to high-level representations associated with meaning in context.

Character Encodings

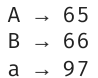
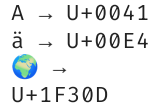
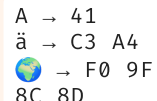
Human writing systems

Societal history of representing meaning with symbols.

Cuneiform c. 3200 BCE 	Hieroglyphs c. 3100 BCE 	Chinese characters c. 1200 BCE - today 
Alphabetic writing c. 700 BCE - today 	Digital-era symbols 20th-21st century 	Learned meaning different fonts, same symbol 

Machine character encodings

Computers encode characters as integer numbers.

ASCII early character numbers 	Unicode code points global character IDs 	UTF-8 bytes stored byte sequences 
--	---	--

What's the result of "A" + "B" in C++? Try it out via [cpp.sh](#)

```
std::cout << ('A' + 'B') << std::endl;
```

Word Tokenizer

Language Models tokenize text

Question: How many "r"s are in Strawberry?

Our view S t r a w b e r r y

Model view [Str] [aw] [berry]

Token IDs [3504, 1134, 19772]

Tokenization depends on the tokenizer.

Modern LLM tokenizers use vocabularies of roughly 100k-200k unique tokens.

Try it yourself: platform.openai.com/tokenizer

Test: Strawberry · Schiffahrt · Poppelsdorf · Poppelsdorfer Allee · She sells seashells.

We also learned to read text in chunks

Hmuans do not raed ervey wrod letetr by letetr. We raed in chnuks, ptatarns, and expcetations. As lnog as the frist and lsat lettres are in the rghit palce, the mnidele can be sracmbeled and the txet is sitll surprisingly easy to raed.

We also use learned representations — not raw pixels or letters alone.

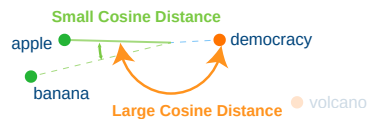
Word Embeddings

From token IDs to vectors

token	ID	vector
strawberry	8123	[0.21, -0.11, ...]
apple	15421	[-0.35, 0.42, ...]
banana	9821	[0.19, 0.05, ...]

Token IDs are arbitrary labels. Embedding vectors are learned representations that can capture meaning.

Words are points in a semantic space

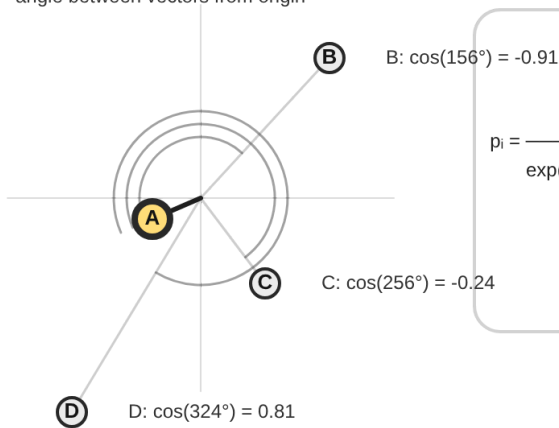


Points are defined by embedding vectors

Similarity can be measured with cosine similarity.

Softmax: Cosine Similarity to Probabilities

Cosine similarity to A
angle between vectors from origin



softmax

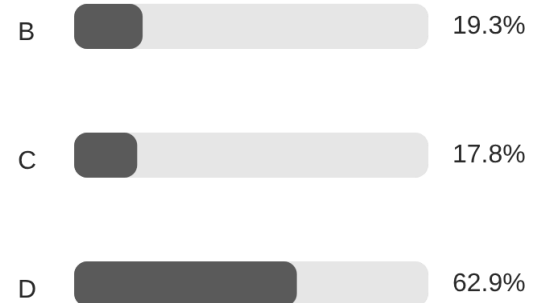
$$p_i = \frac{\exp(\cos^i)}{\exp(\cos^B) + \exp(\cos^C) + \exp(\cos^D)}$$

for $i \in \{B, C, D\}$

$$p^B = \exp(-0.62) / \Sigma = 0.19$$
$$p^C = \exp(-0.71) / \Sigma = 0.18$$
$$p^D = \exp(0.55) / \Sigma = 0.63$$

Probabilities given A

softmax over cosine scores



Drag A, B, C, or D. Probabilities are computed relative to A.

Context matters

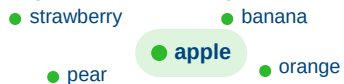
The same word can move to different semantic neighborhoods depending on context.

Fruit context

I ate an **apple** .

apple | "I ate an ..." → fruit-like vector

2D illustration of high-dimensional embeddings



nearby words: fruit, sweet, edible

Technology context

I bought the new **Apple** .

Apple | "I bought the new ..." → company

2D illustration of high-dimensional embeddings



nearby words: device, software, company

Language Models predict the next token.

Static embedding

apple → one vector

Contextual embedding

apple in "I ate an ..." → fruit-like vector
Apple in "I bought the new ..." → company-like vector

Inputs

1. word to predict: <CLS>

2. context: other text

Statistical language model

$p(\langle \text{CLS} \rangle \mid \text{context})$

Output probabilities

fruit: 0.42
company: 0.31
device: 0.18
other: 0.09

Try masked word prediction: huggingface.co/spaces/ysdede/fill-mask-demo

Autoregressive next token prediction

One token at a time

Language models predict the next token from a special masked position and the context that comes before it.

After a token is predicted, it is appended to the sequence. The model then repeats the same prediction step for the next position.

Autoregressive generation means predicting one token after another.

I ate <msk>

Language Model $p(\text{<msk>} \mid \text{I ate})$

softmax over candidate next tokens

1. add mask token
2. estimate probabilities
3. predict mask
4. append to sequence

◀ Previous Token

◀ Previous Step

Next Step ▶

Next Token ▶▶

Training signal

Models are trained with supervised learning to minimize cross-entropy between predicted token probabilities and the actual next token.

Language Model Training

Pre-training

Language models are pre-trained to predict the next tokens as accurately as possible given an existing text corpus.

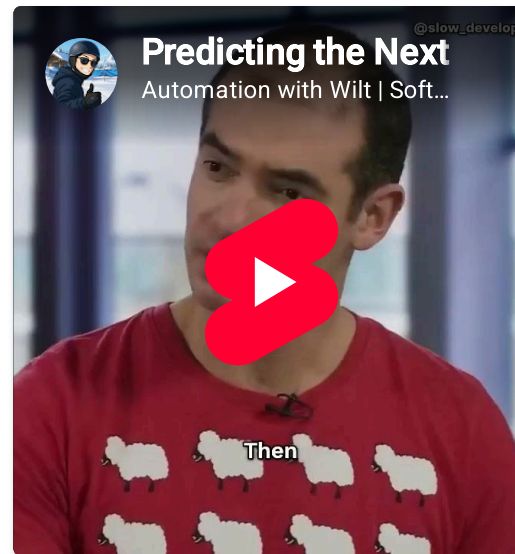
Predicting the next token requires world knowledge.

Fine-tuning

They are then further fine-tuned to match developer-defined prompts.

Fine-tuning creates AI assistants like ChatGPT, Claude, Mistral.

More details: [Andrew Karpathy, “Deep Dive into LLMs like ChatGPT”](#)



[NVIDIA. \(2023\). Fireside Chat With Ilya Sutskever and Jensen Huang: AI Today and Vision of the Future. YouTube short edited by @automationwithwilt.](#)

Geospatial Representations

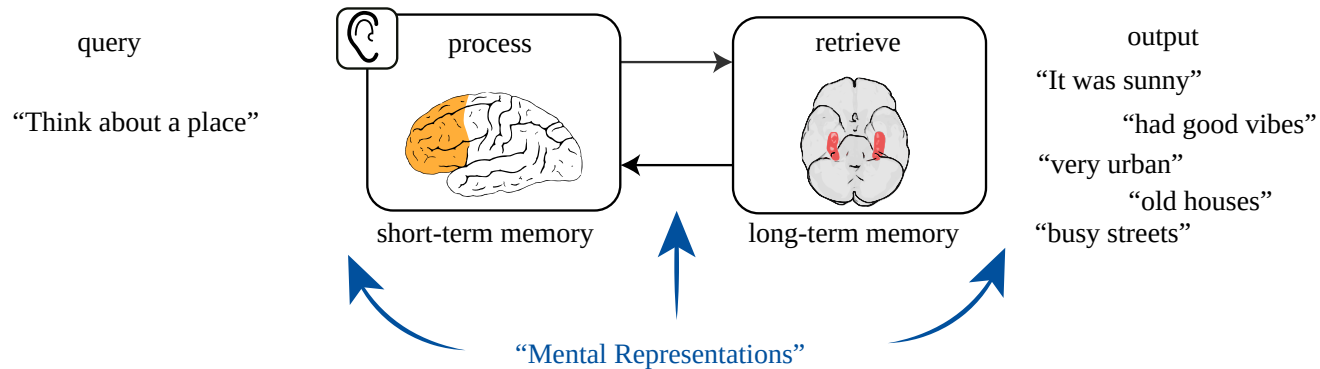


Die Stadt in meinem Kopf. (2016). Süddeutsche Zeitung.

Task: Think of a place

Think of a place you visited on vacation.

Not the country or city name first — think of the actual place. What do you remember?



We feel mental representations as "intuitions" and "memories" of things/places.

A place is more than a coordinate

Coordinate

50.7374° N, 7.0982° E

reference

position

index

Place

home

routes

landmarks

memories

context

meaning

A coordinate tells us where something is. A representation tells us what that place means.

Representations of Geospatial Information

Same city, different maps



We need to select or learn a suitable representation for the problem-at-hand.

How do we experience and recognize places?

We infer location from many weak cues.



vegetation



road markings



architecture



signs and language



terrain



climate



infrastructure

We do not recognize places from one signal. We combine many imperfect signals into a coherent spatial intuition.











Spatial Representations for Earth and Society

Modern geospatial AI aims to learn representations of our planet across space, time, scale, and modality.

Environmental Systems

-  Deforestation and forest degradation
-  Climate change and global warming
-  Flooding and sea level rise
-  Desertification and drought monitoring
-  Glacier retreat and polar ice loss
-  Agriculture and crop monitoring
-  Marine ecosystems and coastal change
-  Environmental pollution and plastic waste
-  Wildfire detection and risk prediction
-  Biodiversity and species distribution

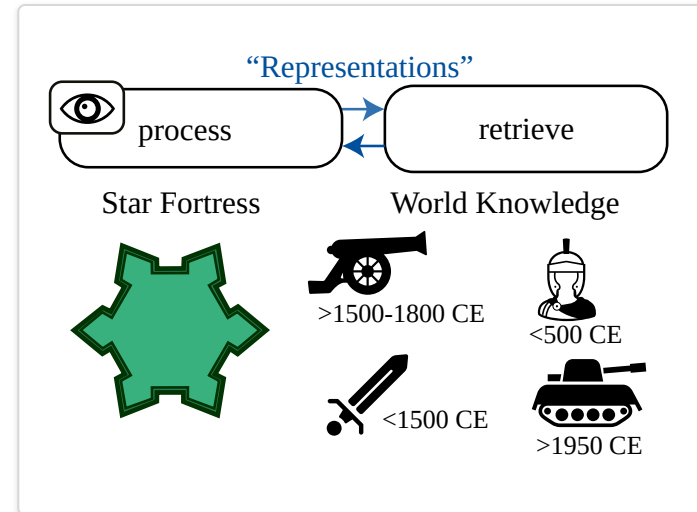
Societal & Human Systems

-  Urbanization and city growth
-  Navigation and route planning
-  Traffic and mobility analysis
-  Disaster response and humanitarian aid
-  Infrastructure and construction monitoring
-  Energy systems and resource management
-  Population density and settlement mapping
-  Migration and human mobility
-  Security and border monitoring
-  Logistics and supply chain optimization

History: When was this city founded?



- 282 CE (Roman/Antique)
- 1153 CE (Medieval)
- 1593 CE (Renaissance)**
- 1975 CE (Modern)



Livability: How livable is this neighborhood?

Livability: How livable is this neighborhood?



- simple
- middle
- good
- very good**

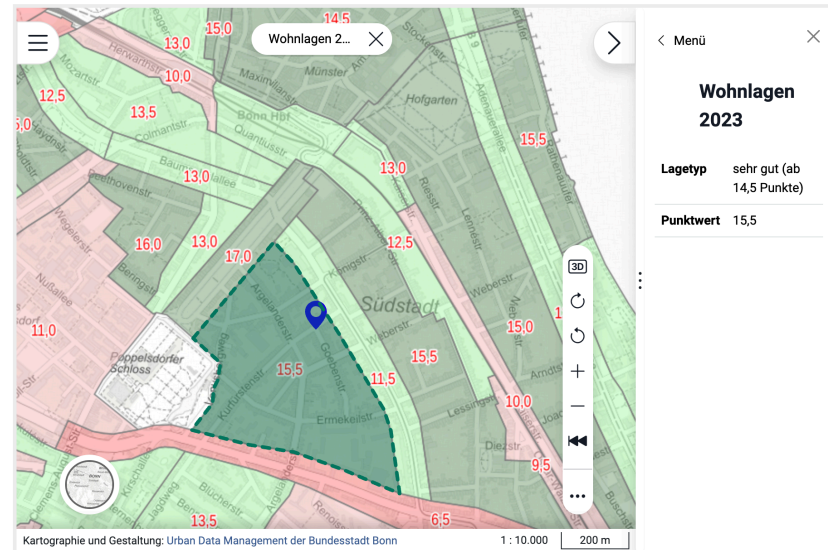
Based on which features in the image, and based on what world knowledge, did you decide?

Livability: Qualitative and Quantitative Factors

The residential location map has four levels: **simple**, **middle**, **good**, and **very good**.

It takes into account:

- Infrastructure
- Urban green space
- **Streetscape**
- Transport connection
- Appreciation
- Burden
- Centrality



Biodiversity: Which Agricultural Landscape supports more Diverse Bird Populations?

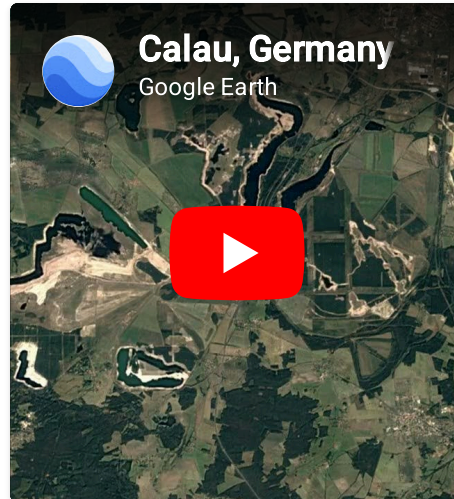
[Engist, D., Finger, R., Knaus, P., Guélat, J., & Wuepper, D. \(2023\). Agricultural systems and biodiversity: evidence from European borders and bird populations. *Ecological Economics*, 209, 107854.](#)

Data of our Spatial Environment

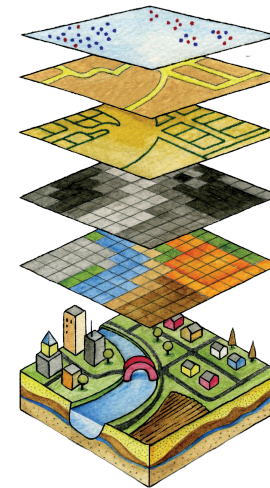
Virtual Reconstructions



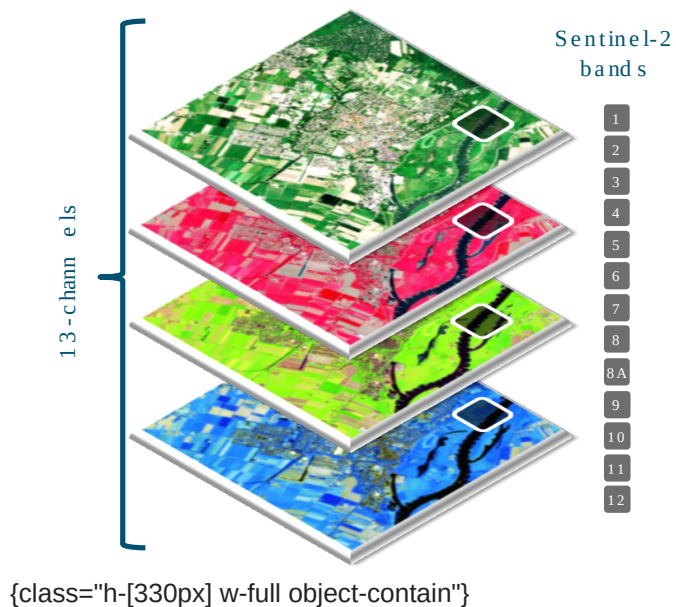
Dense Dynamic Measurements



Geospatial data layers

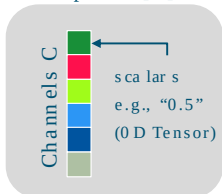


Geodata are Data Tensors

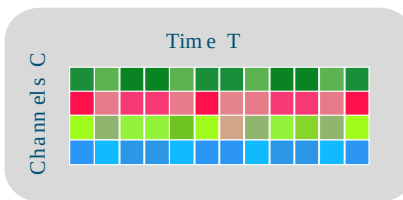


Examples of different data tensors

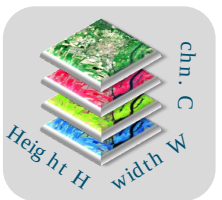
1D Tensor: Vector
e.g., pixels or class probs [C]



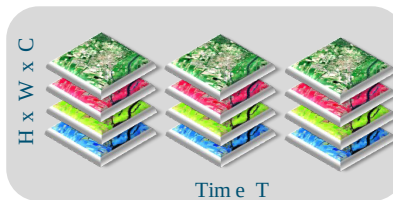
2D Tensor: Matrix
Time Series [T x C]



3D Tensors
Images [H x W x C]



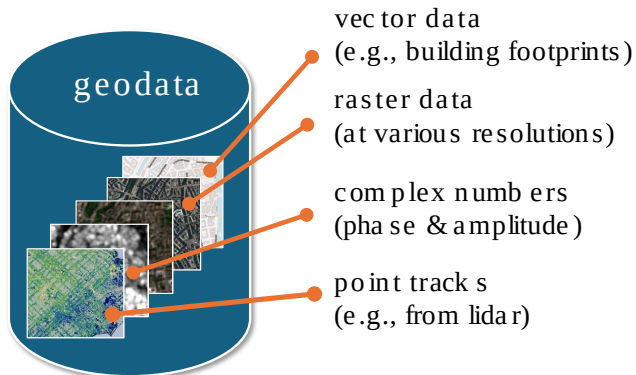
4D Tensors
Img. Time Series [T x H x W x C]



{class="mt-4 h-[300px] w-full object-contain"}

Heterogeneous Datasets in Geodatabases

Geodata is heterogeneous



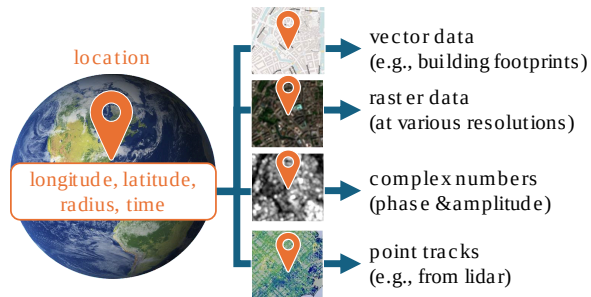
{class="mt-5 h-[330px] w-full object-contain"}

Implementation examples

- **Spatial database:** PostGIS extends PostgreSQL with spatial types, indexes, and functions.
- **Geolocated raster:** GeoTIFF stores georeferenced imagery in TIFF files.
- **Multidimensional arrays:** NetCDF stores scientific data across dimensions such as time, depth, latitude, and longitude.

Different sources use different schemas, coordinate systems, resolutions, timestamps, and file or database abstractions.

Location and Coordinates join all Geodata



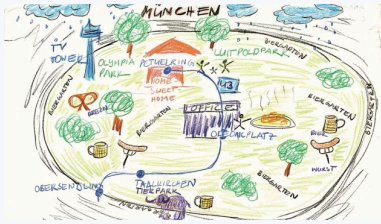
Any place on Earth can be represented by
coordinates

- **Longitude**
- **Latitude**
- **Elevation**
- **Time**

Coordinates join heterogeneous geodata sources together because they reference the same place and time.

Takeaways Mental Representations

Mental representations



Human memories, intuitions, expectations, and place knowledge built from experience.

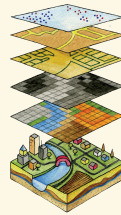
Geospatial representations

Mental representations



Human memories, intuitions, expectations, and place knowledge built from experience.

Geospatial databases



Rasters, vectors, polygons, and geodata layers that encode explicit spatial structure.

Closing and Takeaways

Takeaway: We Navigate the World through Representations



Biological

Perception selects and transforms reality into signals we can act on.



Engineered

Maps, coordinates, images, diagrams, and numbers make selected structure visible.



Learned

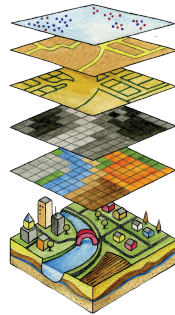
Neural networks learn representations that make patterns useful for prediction.

A representation is a useful reduction of reality: it preserves what matters for a task and hides what does not.

Choosing or Learning Representations

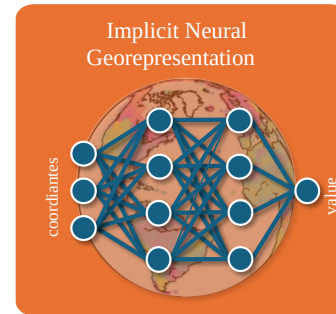
Designed representations

We choose structures that make operations easier.



Learned representations

We train models to discover useful structure.



This course focuses on learning representations.



PRACTICAL 1

EXPLORING GEOSPATIAL REPRESENTATIONS

Geospatial Representation Learning

PRESS – OR SPACE.



© Marc Rußvurm

Licensed under [Creative Commons Attribution–NonCommercial 4.0 International \(CC BY-NC 4.0\)](https://creativecommons.org/licenses/by-nc/4.0/).
You may share and adapt these slides for teaching, research, and other non-commercial educational purposes with attribution.
Commercial training, paid workshops, consulting seminars, or incorporation into commercial course products require prior permission.