

The background of the slide features a stylized landscape. In the foreground, there are dark, rolling hills. In the mid-ground, there are lighter, hazy mountains. On the right side, a large, dark tree with green foliage stands prominently. The sky is white with a few simple, yellow, wavy lines representing clouds. The overall style is minimalist and artistic.

Introduction to Generative AI

Generative AI and Prompt Engineering

Ram N Sangwan



Vectors and Embeddings

Vectors

- On a scale of 0 to 100, how introverted/extraverted are you?
- Have you ever taken a personality test like Big Five Personality Traits ?

These tests ask you a list of questions, then score you on a number of axes, introversion/extraversion being one of them.

Trait	Score
Openness to experience	79 out of 100
Agreeableness	75 out of 100
Conscientiousness	42 out of 100
Negative emotionality	50 out of 100
Extraversion	58 out of 100

Imagine I've scored 38/100 as my introversion/extraversion score. We can plot that in this way.

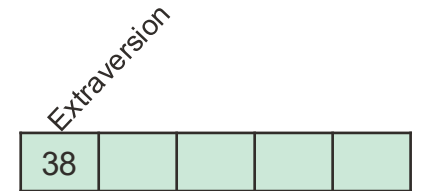
Extraversion

100

0

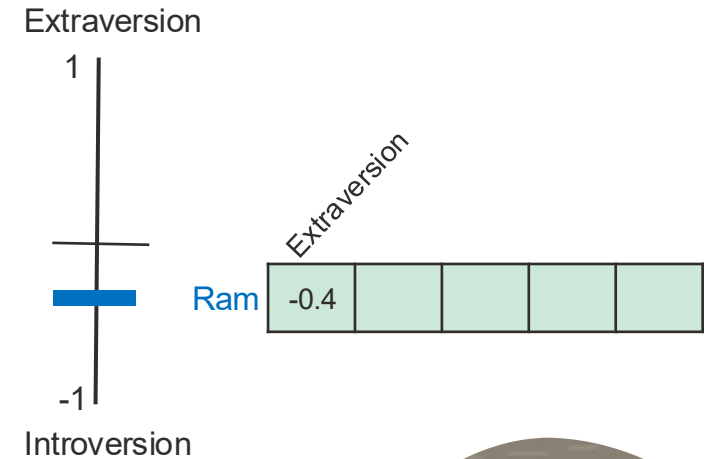
Introversion

Ram



Vectors

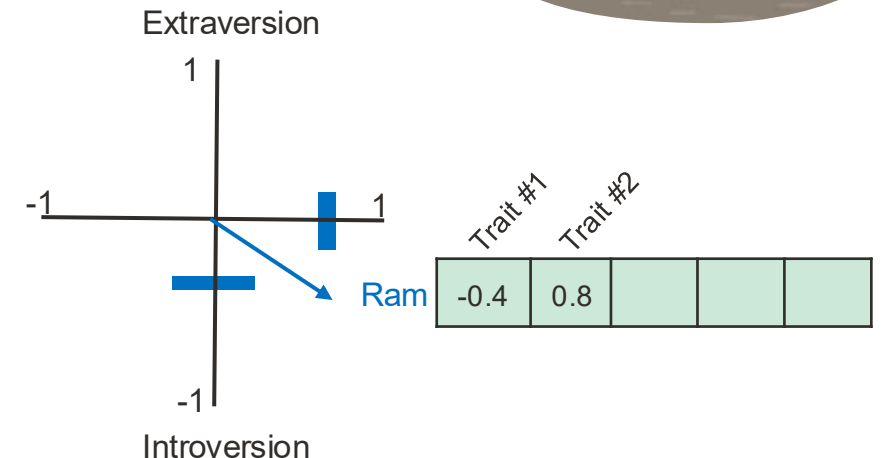
Let's switch the range to be from -1 to 1



- How well do you feel you know a person knowing only this one piece of information?
- Not much.

Let's add another dimension – the score of another trait.

- *We can represent the two dimensions as a vector from the origin to that point.*
- *I've hidden which traits we're plotting so that you get used to not knowing what each dimension represents*
 - *but still getting a lot of value from the vector representation of a person's personality.*



Vectors

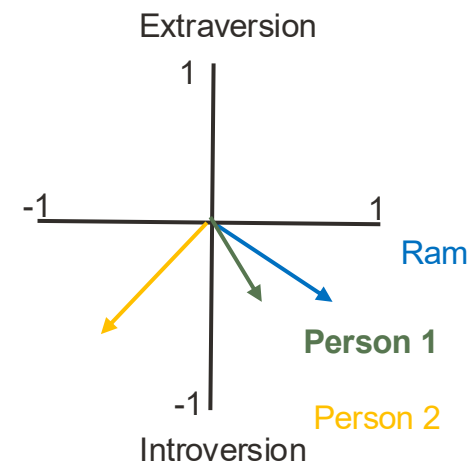
- The usefulness of such representation comes when you want to compare me with others.

A common way to calculate a similarity score for vectors is cosine similarity

Say I am looking for someone with a similar personality.
Which of the two people is more like me?

$$\text{cosine_similarity}(\text{Ram} \begin{bmatrix} -0.4 & 0.8 \end{bmatrix}, \text{Person \#1} \begin{bmatrix} -0.3 & 0.2 \end{bmatrix}) = 0.87 \checkmark$$

$$\text{cosine_similarity}(\text{Ram} \begin{bmatrix} -0.4 & 0.8 \end{bmatrix}, \text{Person \#2} \begin{bmatrix} -0.5 & -0.4 \end{bmatrix}) = -0.20$$



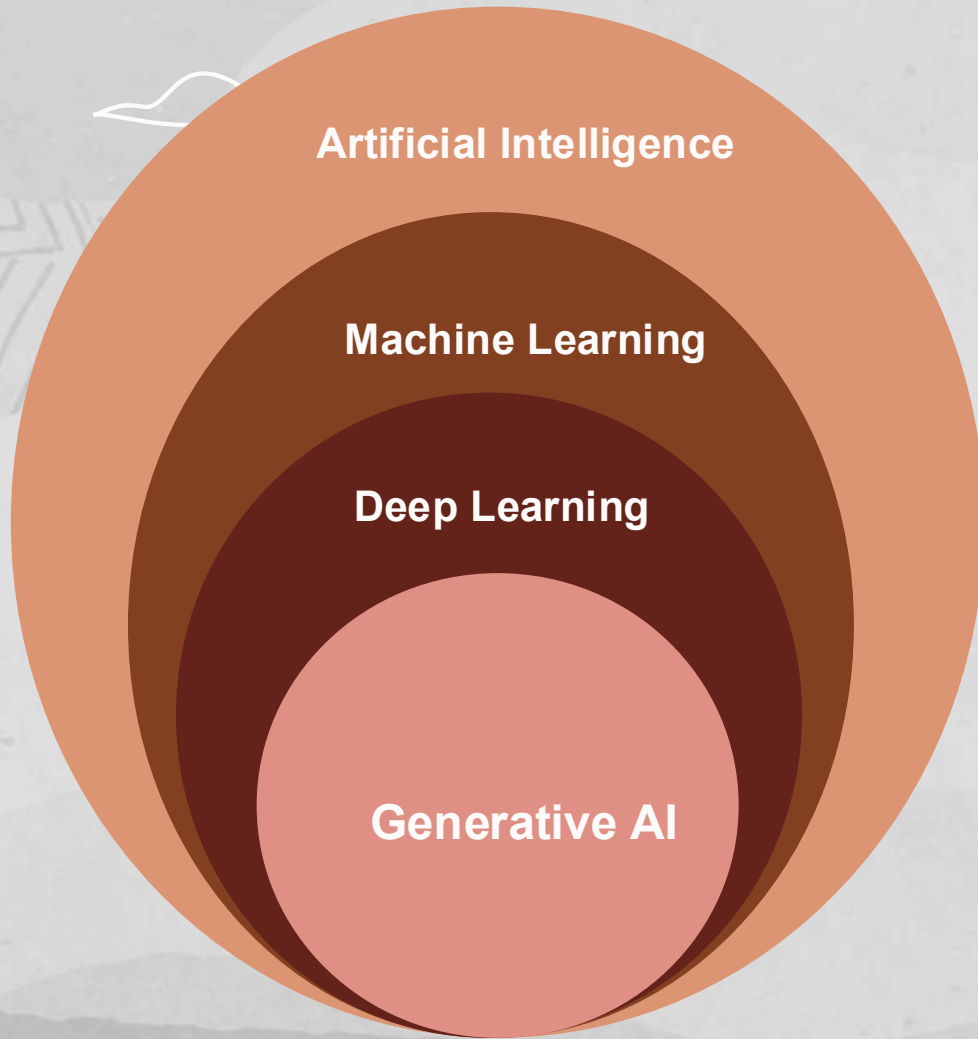
	Trait #1	Trait #2			
Ram	-0.4	0.8			
Person #1	-0.3	0.2			
Person #2	-0.5	-0.4			

- Person #1 is more like me.
- Vectors pointing at the same direction (length plays a role as well) have a higher cosine similarity score.

Now we have two central ideas:

- We can represent people (and things) as vectors of numbers (which is great for machines!).
- We can easily calculate how similar vectors are to each other.

What is Generative AI



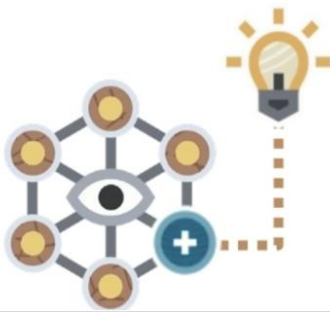
A set of AI methodologies that can create content that resembles the training data they were exposed to.

- A type of AI that can create new content.
- Subset of Deep Learning where the models are trained to generate output on their own.
- Models that can create a wide range of outputs such as images, music, speech, text and other types of data.

How does Generative AI Work?

Learns the underlying patterns in a given data set and uses that knowledge to create new data that shares those patterns.

**Generative
AI Model**



Training Data

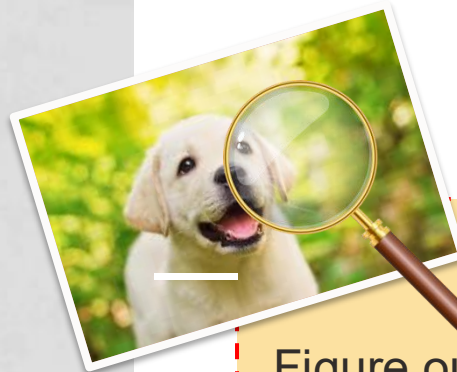
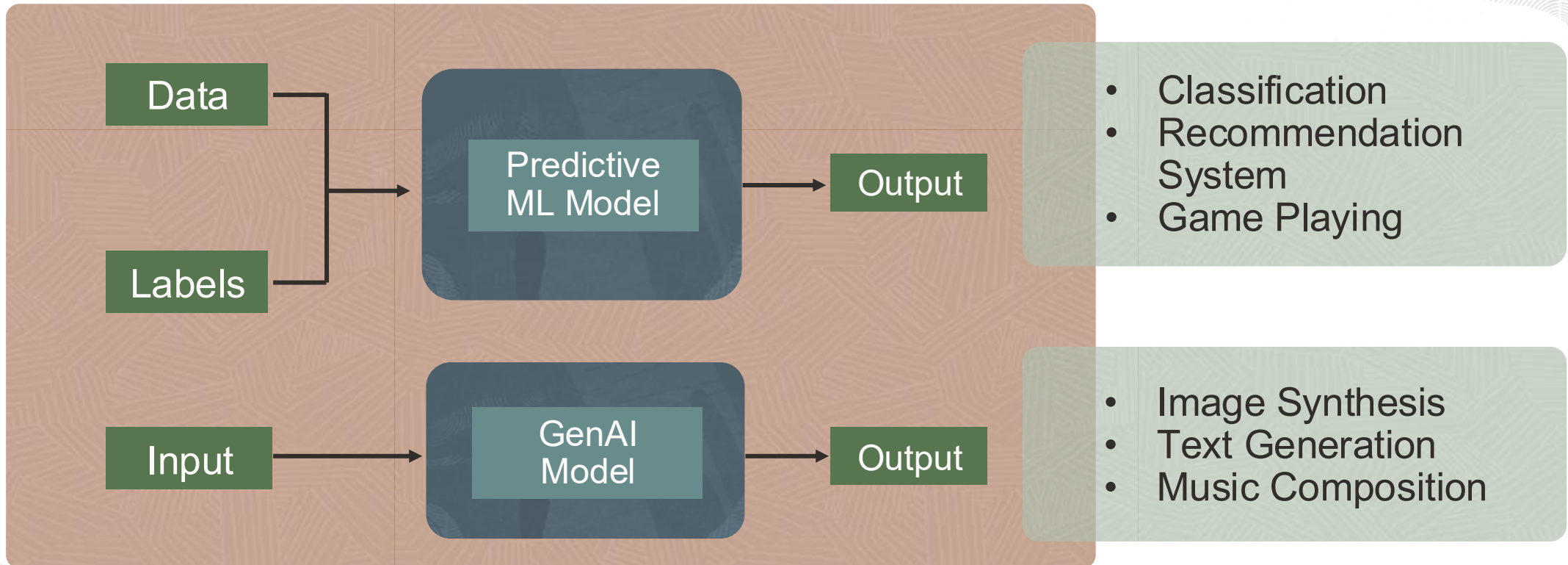


Figure out common dog
patterns and features

“Draw a picture of a Dog”



Generative AI and Other AI Approaches



Types of Generative AI Models



Image Based

- Generates Visual Content
- Learns from Large collection of images

Text-Based

- Generates Textual Content
- Learns from large collection of text data

Generative Adversarial Network (GAN)



- Generate realistic images that resemble training data.
- Create high-quality images and original artworks.
- Imagine you have a bunch of **cat images**, and you want a machine learning model to create similar images.
- This is exactly what a GAN does.
- Price Realised : 432,500 for the image here

<https://www.christies.com/lot/lot-edmond-de-belamy-from-la-famille-de-6166184/>

Generative Adversarial Networks (GANs)

Generator

- Takes in random numbers as input and generates the images of interest (**the forger**)

Discriminator

- Takes both the images from the generator and the real images from the data and spots the difference between them (**the detective**)

Adversarial Objective

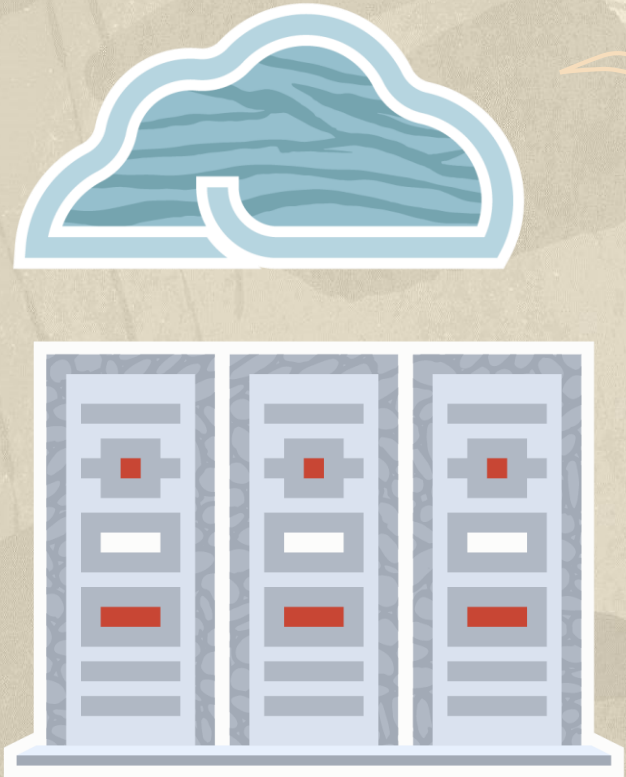
- These two networks are pitted against each other where the **generator creates more realistic synthetic images** to fool the discriminator while the discriminator networks **tries to get better at detecting fake images**.

- Both the generator and the discriminator are trained together.
- And, over the duration of training, the generator gets better at creating images which look real, and the discriminator gets better at spotting fakes.
- This back-and-forth strategy forces both the networks to improve until the generator can create highly realistic synthetic images, that indistinguishable from real images

Diffusion Models



- Work by adding noise to the images in the training data by **forward diffusion process** and then reversing the process to recover the original image using **reverse diffusion**.
- These models can be trained on **large un-labeled datasets** in an **unsupervised manner**.



Generative AI Real-World Use Cases



Visual

- Image generation
- Video Generation
- Design



Language

- Content Creation
- Code Generation
- Conversational AI



Drug Discovery

- New molecular Structure



Music

- Music Generation

Mechanics of Generative AI



Generative Models

Can generate new instances based on what they have learned. E.g. Variational Autoencoders, GANs, and RNNs.

Training Data

The quality of dataset directly impacts the performance of the generated outputs.

Loss Functions

Mathematical functions that measure the difference between the generated output and a desired target.

- Guide the learning process by providing feedback on how well the model is performing.

Optimization Algorithms

Adjust the parameters of the generative model to minimize the loss function during training. E.g. Stochastic Gradient Descent (SGD), Adam, and RMSProp.

Evaluation Metrics

Metrics such as perplexity for language models or Inception Score for image generation tasks.

Hyperparameters and Tuning

Settings that control the behaviour of the learning process. E.g. learning rate, batch size, number of layers in the network, etc.



Mechanics of Generative AI



Regularization Techniques

Help prevent overfitting by adding constraints to the model's parameters or architecture during training. E.g. dropout, weight decay (L2), and early stopping.



Data Augmentation

Involves generating additional training data from existing instances by applying transformations such as rotation, scaling, flipping, etc. It can help improve the generalization ability of generative models.



Thank You

