

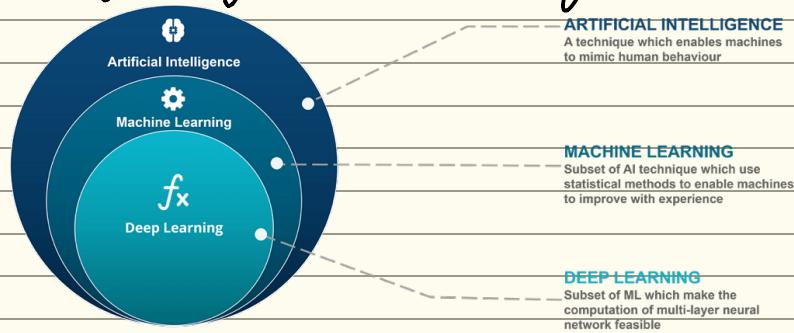
# Azure AI Fundamentals (AI-900)

5 domains that will be tested:

1. Describe AI workloads and Considerations (15-20% 6 to 9 questions)
2. Describe fundamental principles of ML on Azure (20-25% 7-12 questions)
3. Describe features of computer vision workloads on Azure (15-20% 6-9 qs)
4. Describe features of Natural Language Processing (NLP) on Azure (15-20% 6-9 qs)
5. Describe features of generative AI workloads on Azure (15-20% 6-9 qs)

ML and AI Concepts

The layers of Machine learning



Key elements of AI (according to Azure)

- Machine Learning → foundation of an AI sys, learn & predict like a human
- Anomaly Detection → detect outliers or things out of place like a human
- Computer Vision → be able to see like a human
- Natural language processing (NLP) → process human langs & infer context
- Conversational AI → be able to hold conversation with a human

Datasets → logical grouping of units of data

MNIST database → handwritten digits used for classification, clustering & image processing

Common Objects in Context (COCO) dataset → images in a JSON file

\* Data Labelling, Supervised and Unsupervised Learning

Ground Truth → labelled dataset given to model to train and assess.

# Supervised vs Unsupervised vs Reinforcement

Criteria	Supervised ML	Unsupervised ML	Reinforcement ML
Definition	Machine Learns by using labelled data	Machine is trained using unlabelled data without any guidance.	Agent interacts with the environment by performing action. Learns by errors and rewards.
Type of data	Labelled data	Unlabelled data	No - predefined data.
Type of problems	Regression and classification	Association and Clustering	Reward and error based.
Supervision	External supervision	No supervision	No supervision
Algorithms	Linear Regression, Logistic Regression, Naïve Bayes Decision trees	K-Means clustering, KNN (K-nearest neighbours) Principle Component Analysis Neural Networks	Monte Carlo, Q-Learning, SARSA
Aim	Calculate outcomes	Discover underlying patterns	Learn a series of action
Approach	Maps labelled inputs to the known outputs	Understands patterns & discover the output	Follow the trial and error method
Application	Risk Evaluation, Forecast Sales	Recommendation System, Anomaly Detection	Self-Driving Cars, Gaming, Healthcare

Neural Network  $\rightarrow$  Structure that mimicks the brain. It contains layers of nodes called as neurons. The layers are known as perception layers. The connection b/w neurons is weighted

Deep Learning  $\rightarrow$  Neural Network that has 3 or more hidden layers

Feed forward NN  $\rightarrow$  NN where connections b/w nodes don't form a cycle

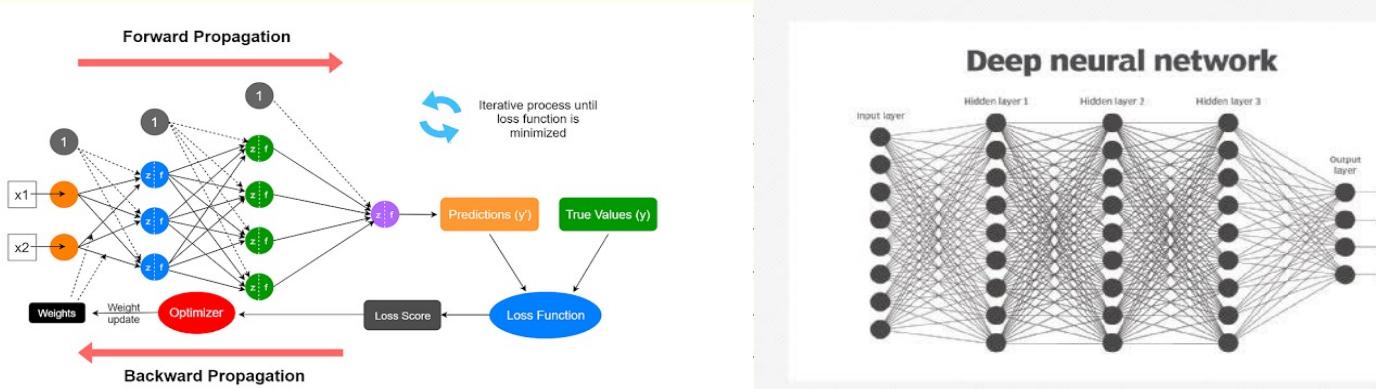
Back propagation  $\rightarrow$  Moves backwards through the NN adjusting weights to improve outcome on next iteration.

Loss function  $\rightarrow$  function that compares ground truth to prediction to determine error rate

Activation functions  $\rightarrow$  Algorithm applied to hidden layer node that affects connected output

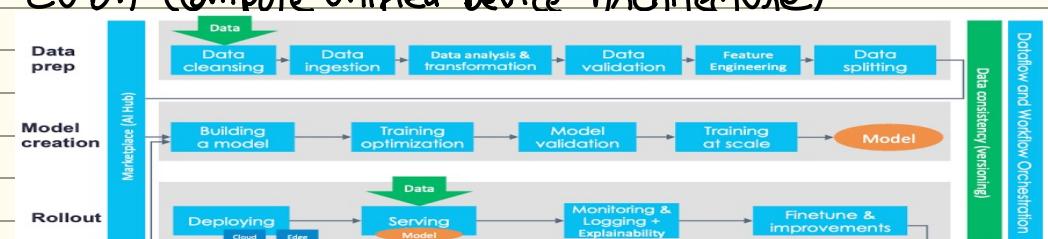
Dense  $\rightarrow$  When the next layer increases the amount of nodes.

Sparse  $\rightarrow$  When the next layer decreases the amount of nodes.



Imp of GPU in ML, CUDA (Compute Unified Device Architecture)

ML Pipeline



# Forecasting vs Prediction

# Talked about JupyterLab and Jupyter Notebooks.

## FORECASTING VERSUS PLANNING

Basis of Comparison	Forecasting	Prediction
Meaning	Process of creating future predictions with relevant data	Process of creating future predictions with or without relevant data
Accuracy	More accurate	Lower probability of happening
Application	Mostly applied in the meteorology, economic and financial sectors	Can be applied almost anywhere
Bias	Forecasts are generated from calculation and data assessment	Is subject to bias
Quantification	Easily Quantified	Can't be quantified
Basis	Done using scientific methods	Arrived at by arbitrary methods e.g. instincts
Application level	Aggregate level	Customer level

DifferenceBetween.net

# Metrics

## What are Metrics?

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

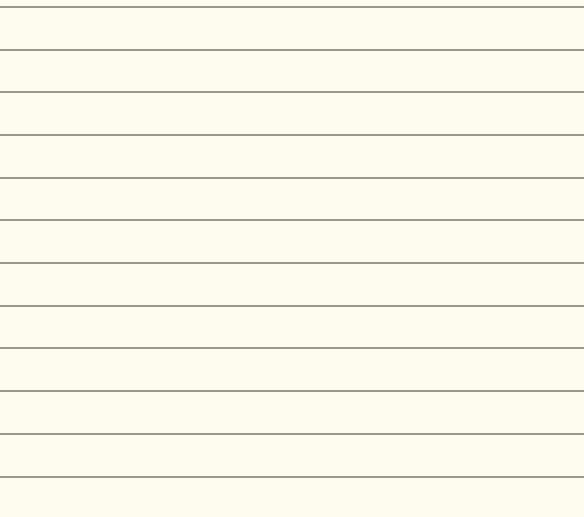
**Performance/Evaluation Metrics** are used to evaluate different Machine Learning Algorithms

For different types of problems different metrics matter, (*this is not an exhaustive list*)

- Classification Metrics (accuracy, precision, recall, F1-score, ROC, AUC)
- Regression Metrics (MSE, RMSE MAE)
- Ranking Metrics (MRR, DCG, NDCG)
- Statistical Metrics (Correlation)
- Computer Vision Metrics (PSNR, SSIM, IoU)
- NLP Metrics (Perplexity, BLEU, METEOR, ROUGE)
- Deep Learning Related Metrics (Inception score, Frechet Inception distance)

There are two categories of evaluation metrics

- Internal Evaluation — metrics used to evaluate the internals of the ML model
  - Accuracy, F1 Score, Precision, Recall (The Famous Four) used in all kinds of models
- External Evaluation — metrics used to evaluate the final prediction of the ML model

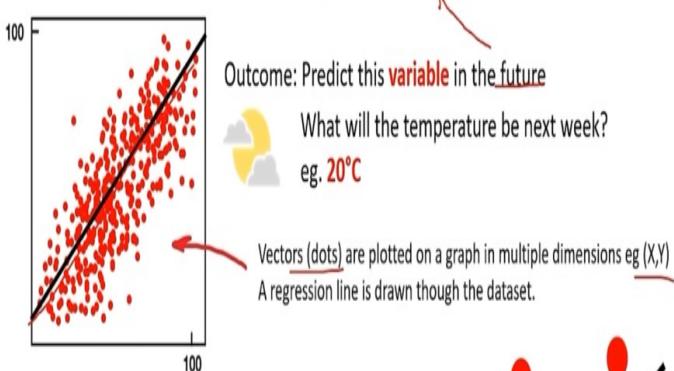


# Regression

## Regression

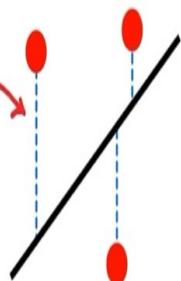
Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Regression is a process of finding a function to correlate a labeled dataset into continuous variable/number.



The distance of the vector from the regression line called an Error  
Different Regression algorithms use the error to predict future variables:

- Mean squared error (MSE)
- Root mean squared error (RMSE)
- Mean absolute error (MAE)



Classification is a process of finding a function to divide a labeled dataset into classes/categories

Outcome: Predict category to apply to the inputted data



### Classification Algorithms

- Logistic Regression
- Decision Tree/Random Forest
- Neural Networks
- Naive Bayes
- K-Nearest Neighbors

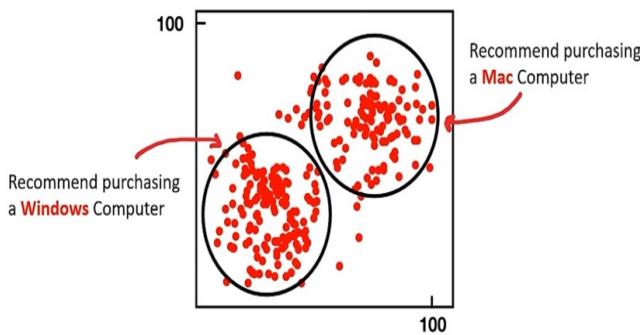
# Clustering

## Clustering

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Clustering is a process grouping unlabeled data based on similarities and differences.

Outcome: Group data based on their similarities or differences



## Confusion Matrix

### Classification Metrics – Confusion Matrix

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

A confusion matrix is a table to visualize the **model predictions** (predicted) vs **ground truth labels** (actual)

Also known as an error matrix. They are useful in classification problems



How many people ate the banana?

	Predicted NO negative	Predicted YES positive	
Actual NO false	75 False Negatives (FN)	25 False Positives (FP)	Our ground truth had 100 labeled items Total False (tF)
Actual YES true	50 True Negatives (TN)	20 True Positives (TP)	Our model made 70 predictions Total True (tT)
100 were NO Total Negative (tN)	75 were YES Total Positive (tP)		

The size of matrix is dependent on the labels:  
Apple, Banana, Orange 3x2 = 6 cells

We have total 170 items  
Total (t)

## Anomaly Detection AI

### Anomaly Detection AI

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

#### What is an anomaly?

An abnormal thing; a marked deviation from the norm or a standard

#### What is anomaly detection?

Anomaly Detection is the process of finding outliers within a dataset called an **anomaly**  
Detecting when a piece of data or access patterns appear suspicious or malicious

Use cases for anomaly detection

- Data cleaning
- **Intrusion detection**
- **Fraud detection**
- Systems health monitoring
- Event detection in sensor networks
- Ecosystem disturbances
- Detection of critical and cascading flaws

Anomaly detection by hand is a very tedious process.  
Using machine learning for anomaly detection is more efficient and accurate



**Anomaly detector** Detect anomalies in data to quickly identify and troubleshoot issues.

# Computer Vision AI

## Computer Vision

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Computer Vision is when we use Machine Learning Neural Networks to gain high-level understanding from digital images or video

### Computer Vision Deep Learning Algorithms:

- Convolutional neural network (CNN) — image and video recognition
  - Inspired after how human eyes actually process information and send it back to brain to be processed
- Recurrent neural network (RNN) — handwriting recognition or speech recognition

### Types of Computer Vision

- Image Classification — look at an image or video and classify (place it in a category)
- Object Detection — identify objects within an image or video and apply labels and location boundaries
- Semantic Segmentation — identify segments or objects by drawing pixel mask (great for objects in movement)
- Image Analysis — analyze an image or video to apply descriptive and context labels
  - eg. An employee sitting at a desk in Tokyo
- Optical Character Recognition — Find text in images or videos and extract them into digital text for editing
- Facial Detection — detect faces in a photo or video, draw a location boundary, label their expression



Seeing AI is an AI app developed by Microsoft for iOS

Seeing AI uses the device camera to **identify people and objects**, and then the app audibly **describes those objects for people with visual impairment**.

## Computer Vision

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

### Azure's Computer Vision Service Offering:



Computer Vision analyze images and video, and extract descriptions, tags, objects, and text



Custom Vision custom image classification and object detection models using your own images



Face Detect and identify people and emotions in images.



Form Recogniser translate scanned documents into key / value or tabular editable data

# Natural Language Processing (NLP)

## Natural Language Processing (NLP)

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Natural Language Processing is Machine Learning that can **understand the context of a corpus (a body of related text)**.

NLP enables you to:

- Analyze and interpret text within documents, email messages
- Interpret or contextualise spoken token eg sentiment analysis
- Synthesize speech eg. a voice assistance talking to you
- Automatically translate spoken or written phrases and sentences between languages.
- Interpret spoken or written commands and determine appropriate actions.



Hi, I'm Cortana.

Cortana is a **virtual assistant** developed by Microsoft which uses the Bing search engine to perform tasks such as setting reminders and answering questions for the user.

## Natural Language Processing (NLP)

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

### Azure's NLP Service Offering:



#### Text Analytics

- sentiment analysis to find out what customers think
- Find topic-relevant phrases using key phrase extraction
- identify the language of the text with language detection
- Detect and categorize entities in your text with named entity recognition



#### Translator

- real-time text translation
- multi-language support



#### Speech

- transcribe audible speech into readable, searchable text



#### Language Understanding (LUIS)

- natural language processing service that enables you to understand human language in your own application, website, chatbot, IoT device, and more

# Conversational AI

## Conversational AI

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Conversational AI is technology that can **participate in conversations with humans**.

- Chatbots
- Voice Assistants
- Interactive Voice Recognition Systems (IVRS)

### Use Cases

- Online Customer Support — replaces human agents for replying about customer FAQs, shipping
- Accessibility — voice operated UI for those who are visually impaired
- HR processes — employee training, onboarding, updating employee information
- Health Care — accessible and affordable health care eg. claim processes
- Internet of Things (IoT) — Amazon Alexa, Apple Siri and Google Home
- Computer Software — autocomplete search on phone or desktop



QnA Maker Create a conversational question-and-answer bot from your existing content (knowledge base).



Azure Bot Service Intelligent, serverless bot service that scales on demand. Used for creating, publishing, and managing bots

# Responsible AI

## Responsible AI

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Responsible AI focuses on **ethical, transparent and accountable** use of AI technologies

Microsoft puts into practice Responsible AI via its six **Microsoft AI principles**

1. Fairness — AI systems should treat all people fairly
2. Reliability and Safety — AI systems should perform reliably and safely
3. Privacy and Security — AI systems should be secure and respect privacy
4. Inclusiveness — AI systems should empower everyone and engage people
5. Transparency — AI systems should be understandable
6. Accountability — People should be accountable for AI systems

### Fairness

#### Responsible AI – Fairness

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

##### AI systems should treat all people fairly

AI systems can reinforce existing societal stereotypical  
Bias can be introduced during the development of a pipeline

- AI systems that are used to allocate or withhold:
- opportunities
  - resources
  - Information
  - In domains:
    - Criminal Justice
    - Employment and Hiring
    - Finance and Credit

eg. an ML model designed to select final applicants for a hiring pipeline without incorporating any bias based on gender, ethnicity or may result in an unfair advantage

Azure ML can tell you how each feature can influence a model's prediction for bias

**= Fairlearn** Fairlearn is an open-source python project to help data scientist to improve fairness in their AI systems

### Reliability and Safety

#### Responsible AI – Reliability and safety

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

##### AI systems should perform reliably and safely

AI software must be **rigorous tested** to ensure they work as expected before release to the end user

If there are scenarios where AI is making mistakes its important to release a report **quantified risks and harms** to end-users so they are informed of the short-comings of an AI solution

AI where concern for reliability and safety for humans is critically important:

- Autonomous Vehicle
- AI health diagnosis, AI suggesting prescriptions
- **Autonomous Weapon Systems**

### Privacy and Security

#### Responsible AI – Privacy and security

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

##### AI systems should be secure and respect privacy

AI can require vast amounts of data to train Deep Learning ML models.  
The nature of the ML model may require **Personally identifiable information (PII)**

It is important that we ensure protection of user data that it is not leaked or disclosed

In some cases ML Models can be run locally on a user's device so their PII remains on their device avoiding that vulnerability

- AI Security Principles to detect malicious actors:
- Data Origin and Lineage
  - Data Use Internal vs External
  - Data Corruption Considerations
  - Anomaly detection

### Inclusiveness

#### Responsible AI – Inclusiveness

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

##### AI systems should empower everyone and engage people

##### Minority Groups

If we can design AI solutions for the **minority** of users Then we can design AI solutions for the majority of users

- physical ability
- gender
- sexual orientation
- ethnicity
- other factors

### Transparency

#### Responsible AI – Transparency

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

##### AI systems should be understandable

Interpretability / Intelligibility is when end-users can understand the behaviour of the UI

Transparency of AI systems can result in

- Mitigating unfairness
- Help developers debug their AI systems
- Gaining more trust from our users

Those build AI systems should be:

- open about the why they are using AI
- open about the limitations of their AI systems

Adopting an open-source AI framework can provide transparency (at least from a technical perspective) on the internal workings of an AI systems

### Accountability

#### Responsible AI – Accountability

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

##### People should be accountable for AI systems

The structure put in place to consistently enacting AI principles and taking them into account

AI systems should work within:

- framework of governance
- organizational principles

ethical and legal standards that are clearly defined

Principles guide Microsoft on how they **Develop, Sell and Advocate** when working with third-parties and this can push towards regulations towards AI Principles

Guidelines for Human AI interaction → Microsoft has a website for it.

## Azure Cognitive Services



Azure Cognitive Services is a **comprehensive family of AI services** and cognitive APIs to help you build intelligent apps

- Create customizable, pretrained models built with "breakthrough" AI research
- Deploy Cognitive Services anywhere from the cloud to the edge with containers
- Get started quickly—no machine-learning expertise required
- Developed with strict ethical standards, empowering responsible use with industry-leading tools and guidelines



### Decision

- Anomaly Detector — Identify potential problems early on.
- Content Moderator — Detect potentially offensive or unwanted content.
- Personaliser — Create rich, personalised experiences for every user.

### Language

- Language Understanding — Build natural language understanding into apps, bots and IoT devices.
- QnA Maker — Create a conversational question and answer layer over your data.
- Text Analytics — Detect sentiment, key phrases and named entities.
- Translator — Detect and translate more than 90 supported languages.

### Speech

- Speech to Text — Transcribe audible speech into readable, searchable text.
- Text to Speech — Convert text to lifelike speech for more natural interfaces.
- Speech Translation — Integrate real-time speech translation into your apps.
- Speaker Recognition — Identify and verify the people speaking based on audio.

### Vision

- Computer Vision — Analyze content in images and video.
- Custom Vision — Customize image recognition to fit your business needs.



# Azure Cognitive Services

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Cognitive Services is an umbrella AI service that enables customers to **access multiple AI services** with an **API key and an API Endpoint**

Home >

## Cognitive Services

Exapro Training Inc (exampro.onmicrosoft.com)

+ Add Manage view Refresh Export to CSV

Filter for any field... Subscription == all Resource group

Showing 1 to 1 of 1 records.

Name ↑
<input type="checkbox"/> myCognitiveServices734

Show Keys

KEY 1  
.....

KEY 2  
.....

Endpoint

Location



# Face Service

## Face Service

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Azure Face service provides AI algorithms that **detect, recognize, and analyze human faces** in images

Azure Face can detect:

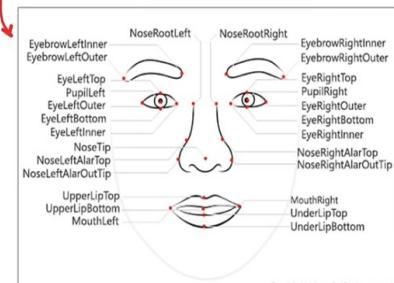
- faces in an image
- faces with specific attributes
- face landmarks
- similar faces
- the same face as a specific identity across a gallery of images



**Face ID**  
unique identifier string for each detected face in an image

### Face Landmarks

easy-to-find points on a face  
27 predefined landmark points.

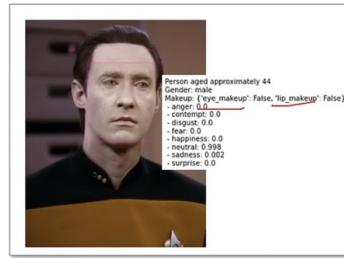


## Face Service

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

### Face Attributes

- Accessories. (Wearing accessories)
- Age
- Blur (blurriness of the face in the image)
- Emotion.
- Exposure
- Facial hair
- Gender
- Glasses
- Hair
- Head pose
- Makeup
- Mask. (are they wearing a mask?)
- Noise. The visual noise detected in the face image
- Occlusion. (objects blocking parts of the face)



# Speech and Translate Service

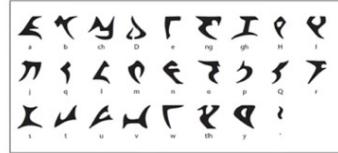
## Speech and Translate Service

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)



Azure's Translate service is a **translation service**.

- It can translate 90 languages and dialects
  - It even supports **Klingon!**
- It uses Neural Machine Translation (NMT) replacing its legacy Statistical Machine Translation (SMT)
- Custom Translator** allows you to extend the service for translation based on your business and domain use case



Azure Speech service can **speech synthesis service** speech-to-text, text-to-speech, and speech-translation

### Speech-to-Text

- Real-time Speech-to-text
- Batch Speech-to-Text
- Multi-device Conversation
- Conversation Transcription
- Create Custom Speech Models

### Text-to-Speech

- using Speech Synthesis Markup Language (SSML)
- Create Custom Voices

### Voice Assistance

### Speech Recognition

- integrates with Bot Framework SDK
- Speaker verification & identification

# Text Analytics

## Text Analytics

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Text Analytics API is a **Natural Language Processing (NLP)** service for **text mining and text analysis**

Text Analytics can perform:

- sentiment analysis**
  - find out what people think of your brand or topic
    - feature provides sentiment labels (such as "negative", "neutral" and "positive")
- opinion mining**
  - aspect-based sentiment analysis
  - granular information about the opinions related to aspects
- key phrase extraction**
  - quickly identify the main concepts in text.
- language detection**
  - detect the language an input text is written in
- named entity recognition (NER)**
  - Identify and categorize entities in your text as people, places, organizations, quantities
  - Subset of NER is Personally Identifiable Information (PII)

# More on NLP services

## NLP – Key Phrase Extraction

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

### Key Phrase Extraction quickly identify the main concepts in text

- Key phrase extraction works best when you give it bigger amounts of text to work on
- This is opposite from sentiment analysis, which performs better on smaller amounts of text
- Document size must be 5,120 or fewer characters per document, and you can have up to 1,000 items (IDs) per collection

When the Borg launch an attack on Earth, the Enterprise is sent to the neutral zone due to the Admiralty's mistrust of Picard's abilities as he had been assimilated in the past. The Enterprise however, disobeys and returns to help destroy the Borg ship. However a smaller ship escapes and travels back in time, causing the assimilation of Earth in the future. The Enterprise follows the ship back in time and have to undo the damage the ship did on the surface to an experimental warp drive unit that will lead Earth to its first contact with alien life. Meanwhile, on the Enterprise, survivors of the Borg ship begin to assimilate decks within the ship itself...

Key Phrases:
Borg ship
Enterprise
smaller ship escapes
time
assimilation of Earth
surface
experimental warp drive unit
Admiralty's mistrust of Picard's abilities
neutral zone
travels
contact
damage
attack
survivors
decks
alien life
future
past

## Named Entity Recognition

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Named Entity Recognition detects **words and phrases mentioned in unstructured text** that can be **associated with one or more semantic types**.

Ribavirin [UMLS: C0035525] was also evaluated against SARS-CoV-2 infection, but the antiviral [UMLS: C0003451] MEDICATION\_NAME MEDICATION\_CLASS

property of drugs [UMLS: C0013227] is still not well established against the SARS-CoV-2 [UMLS: CS203670] negation TREATMENT\_NAME DIAGNOSIS

In addition, after oral administration, the drug was rapidly absorbed into the GI tract [UMLS: C0017189] BODY\_STRUCTURE

The drug has oral bioavailability around 64% with large volume of distribution.

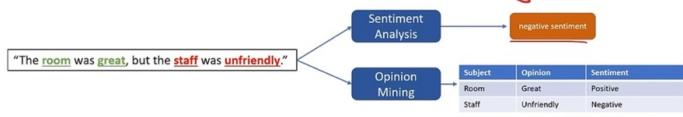
Semantic types could be: Location, Event, Location, Person, Diagnosis, Age

## NLP – Sentiment Analysis

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Sentiment analysis will apply labels and confidence score to text at the **sentence and document level**.

- Labels include **negative, positive, mixed or neutral**
- Confidence scores ranging from 0 to 1



Opinion mining will provide more granular data with a **Subject** and **Opinion** tied to a Sentient

# Optical Character Recognition (OCR)

## Optical Character Recognition (OCR)

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Optical character recognition (OCR) is the process of **extracting printed or handwritten text into a digital and editable format**

- OCR can be applied to:
- photos of street signs
  - **Products** →
  - Documents
  - Invoices
  - Bills
  - Financial Reports
  - Articles
  - and more



## Optical Character Recognition (OCR)

Cheat sheets, Practice Exams and Flash cards  [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Azure has **two different APIs** that can perform OCR: **OCR API** and **Read API**

### OCR API

- older recognition model
- supports only images
- executes synchronously
  - returning immediately with the detected text
  - Suited for less text
- Support more languages
- Easier to implement

### Read API

- updated recognition model
- Supports images and PDFs
- Executes asynchronously
  - parallelizes tasks per line for faster results
  - Suited for lots of text
- Supports fewer languages
- A bit more difficult to implement



OCR is performed via the Computer Vision SDK



# Luis

## Language Understanding Service (Luis)

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Language Understanding (Luis)** is a no-code ML service to build natural language into apps, bots, and IoT devices.

Quickly create enterprise-ready, custom models that continuously improve.

Luis is accessed via its own isolate domain at [luis.ai](https://luis.ai)

Luis utilizes Natural Language Processing (NLP) and **Natural Language Understanding (NLU)**

NLU is the ability to *transform* a linguistic statement to a representation that enables you to understand your users naturally

Luis is intended to focus on **intention** and **extraction**:

- What the user wants
- What they are talking about

## Language Understanding Service (Luis)

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

A Luis application is composed of a **schema**

This schema is autogenerated for you when you use the Luis.ai web interface

The schema defines:

- **Intentions** — what the user is asking for
  - A Luis app always contains a **None** Intent
- **entities** — what parts of the intent is used to determine the answer
- **utterances** — Examples of user input that includes intent and entities to train the ML model to match predictions against real user input
  - An intent requires one or more example utterance for training
    - It is recommended to have 15-30 example utterances
  - To explicitly train to ignore an utterance use the **None** Intent

Intents **classify** user utterances  
Entities **extract** data from utterance



Example Utterance: book me two flights to Toronto

Intent: bookFlight

Entities: two, flights, to, Toronto

# QnA Maker Service

## QnA Maker Service

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**QnA Maker** is a cloud-based Natural Language Processing (NLP) service that allows you to **create a natural conversational layer** over your data.

QnAMaker is hosted on its own isolate domain at [www.qnamaker.ai](https://www.qnamaker.ai)

It will find the most appropriate answer for any input from your **custom knowledge base** (KB) of information

Commonly used to build conversational client applications, which include:

- social media applications
- chat bots
- speech-enabled desktop applications

QnA Maker doesn't store customer data  
All customer data is stored in the region the customer deploys the dependent service instances in

## QnA Maker Service – Knowledgebase

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

QnA Maker imports your content into a knowledge base of question and answer pairs.

QnA Maker can build you knowledge base from an **existing document, manual or website (URL, DOCX, PDF)**

It will use ML to extract the question and answer pairs.

The content of the question and answer pair includes:

- All the alternate forms of the question
- Metadata tags used to filter answer choices during the search
- Follow-up prompts to continue the search refinement

QnA Maker stores answer text as **markdown**



## QnA Maker Service – Chat box

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

You converse with your bot through a Chat Box. There are many opportunities to interact with your bot in QNAMaker.ai, Azure Bot Service, Bot Composer.

Via Channels you can even get embeddable chatbox code



## QnA Maker Service – Layered Ranking

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

QnA Maker's system is a layered ranking approach.

The data is stored in Azure search, which also serves as the first ranking layer.

The top results from Azure search are then passed through QnA Maker's NLP re-ranking model to produce the final results and confidence score.

## QnA Maker Service – Knowledgebase

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Once your Knowledge Base is imported you can **fine-tune the imported results** by editing the Question and Answer pairs

Question	Answer	Metadata tags
I accidentally deleted a part of my QnA Maker, what should I do?	All deletes are permanent, including question and answer pairs, files, URLs, custom questions and answers, knowledge bases, or Azure resources. Make sure you export your knowledge base from the **Settings** page before deleting any part of your knowledge base.	Type : troubleshooting Format : text-only Nextstep : recover
Can I undo deleted questions and answers?		

## QnA Maker Service – Chit Chat

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Chit-chat

- None
- Professional
- Friendly
- Witty
- Caring
- Enthusiastic

The chit-chat feature in QnA maker allows you to easily add a **pre-populated set of the top chit-chat**, into your knowledge base.

This dataset has about **100 scenarios** of chit-chat in the voice of multiple personas

## QnA Maker Service – Multi-turn conversation

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Multi-turn conversation** is follow-up prompts and context to manage the multiple turns, known as *multi-turn*, for your bot from one question to another

When a question **can't be answered in a single turn**

QnA Maker provides multi-turn prompts and active learning to help you improve your basic question and answer pairs.

**Multi-turn prompts** give you the opportunity to connect question and answer pairs. This connection allows the client application to provide a top answer and provides more questions to refine the search for a final answer.

After the knowledge base receives questions from users at the published endpoint, QnA Maker applies **active learning** to these real-world questions to suggest changes to your knowledge base to improve the quality.

# Azure Bot Service



## Azure Bot Service

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)



**Azure Bot Service** Intelligent, serverless bot service that scales on demand.  
Used for creating, publishing, and managing bots

You can register and publish a variety of bots from the Azure Portal

Azure Bot Service can integrate your bot with other Azure, Microsoft or Third Party services via **Channels**:

- Direct Line
- Alexa
- Office 365 email
- Facebook
- Kik
- LINE
- Microsoft Teams
- Skype
- Twilio
- and more....

## Bot Framework Composer

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)



**Bot Framework Composer**, built on the Bot Framework SDK, is an **open-source IDE for developers** to author, test, provision and manage conversational experiences.

Composer is **downloadable app** available for Windows, OSX and Linux

- You can use either C# or Node to build your bot
- Deploy your bots to:
  - Azure Web App
  - Azure Functions
- Templates to build:
  - QnA Maker Bot
  - Enterprise or Personal Assistant Bot
  - Language Bot
  - Calendar or People Bot
- Test and debug via the Bot Framework Emulator
- Built in Package manager

## Bot Framework SDK

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

The Bot Framework SDK v4 is an **open-source SDK** that enables developers to **model and build sophisticated conversations**

The Bot Framework, along with the Azure Bot Service, provides an **end-to-end workflow**:



With this framework, developers can create bots that use speech, understand natural language, handle questions and answers, and more.

The Bot Framework includes a modular and extensible SDK for building bots, as well as tools, templates, and related AI services.

# Azure Machine Learning Service



## Azure Machine Learning Service

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

### Azure Machine Learning Studio (classic)

An older service that manages AI/ML workloads. Does not have a pipeline and other limitations. Workloads are not easily transferable from classic to the new service.

### Azure Machine Learning Service

A service that simplifies running AI/ML related workloads allowing you to build flexible Automated ML Pipelines. Use Python or R, Run DL workloads such as Tensorflow

#### Jupyter Notebooks

- build and document your machine learning models as you build them, share and collaborate

#### Azure Machine Learning SDK for Python

- An SDK designed specifically to interact with Azure Machine Learning Services

#### MLOps

- end to end automation of ML model pipelines eg. CI/CD, training, inference

#### Azure Machine Learning Designer

- drag and drop interface to visually build, test, and deploy machine learning models

#### Data Labeling Service

- ensemble a team of humans to label your training data

#### Responsible Machine Learning

- model fairness through disparity metrics and mitigate unfairness

## Azure Machine Learning Studio – Compute

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Azure Machine Learning Studio has **4 kinds of compute**:

1. **Compute Instances** — Development workstations that data scientists can use to work with data and models.
2. **Compute Clusters** — Scalable clusters of virtual machines for on-demand processing of experiment code.
3. **Inference Clusters** — Deployment targets for predictive services that use your trained models.
4. **Attached Compute** — Links to existing Azure compute resources, such as Virtual Machines or Azure Databricks clusters.

## Azure Machine Learning Studio – Overview

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

### Author

**Notebooks** — Jupyter Notebooks, an IDE to write python code to build ML models

**AutoML** — Completely automated process to build and train an ML model

**Designer** — Visual drag and drop designer to construct end to end ML pipelines

### Assets

**Datasets** — data that you upload which will be used for training

**Experiments** — when you run a training job they are detailed here

**Pipelines** — ML workflows you have built, or you have used in the Designer

**Models** — a model registry containing trained models that can be deployed

**Endpoints** — when you deploy a model its hosted on an accessible endpoint eg. REST API

### Manage

**Compute** — the underlying computing instances used to for notebooks, training, inference

**Environments** — a reproducible Python environment for machine learning experiments

**Datastores** — a data repository where your dataset resides

**Data Labeling** — have humans with ML-assisted labeling to label your data for supervised learning

**Linked Services** — external services you can connect to the workspace eg. Azure Synapse Analytics

## Azure Machine Learning Studio – Data Labeling

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Create Data labeling jobs to prepare your Ground Truth for supervised learning

### Human-in-the-loop labeling

You have a team of humans that will apply labeling  
These are humans you grant access to labeling

### Machine-learning-assisted data labeling

You will use ML to perform labeling

You can export the label data for Machine Learning experimentation at any time  
Users often export multiple times and train different models, rather than wait for all the images to be labeled.

Image labels can be exported in:

- COCO format
- Azure Machine Learning dataset
  - dataset format makes it easy to use for training in Azure Machine Learning

# Azure Machine Learning Studio – Data Stores

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

**Datastores securely connect to your storage service** on Azure without **putting your authentication credentials** and the integrity of your original data source **at risk**.

Datastore type *
Azure Blob Storage
Azure Blob Storage
Azure file share
Azure Data Lake Storage Gen1
Azure Data Lake Storage Gen2
Azure SQL database
Azure PostgreSQL database
Azure MySQL database

- Azure Blob Storage**  
data is stored as objects, distributed across many machines
- Azure File Share**  
a mountable file share via SMB and NFS protocols
- Azure Data Lake Storage (Gen 2)**  
Azure Blob storage designed for vasts amount of data for Big Data analytics
- Azure SQL database**  
Full-managed MS SQL relational database
- Azure Postgres database**  
open-source relational database
- Azure MySQL Database**  
Open-source relational database

# Azure Machine Learning Studio – Datasets

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Azure provides a same code snippet with the **Azure Machine Learning SDK for Python** to start programmatically using datasets in your Jupyter Notebooks

```
Sample usage □  
# azureml-core of version 1.0.72 or higher is required  
# azureml-datatools of version 1.1.34 or higher is required from azureml.core import Workspace, Dataset  
  
subscription_id = '7f0352cf-6c7d-456a-8ec0-83ef2120997b'  
resource_group = 'MyStudio'  
workspace_name = 'MyStudio'  
  
workspace = Workspace(subscription_id, resource_group, workspace_name)  
dataset = Dataset.get_by_name(workspace, name='Sample: Diabetes')  
dataset.to_pandas_dataframe()
```

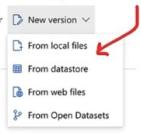
# Azure Machine Learning Studio – Datasets

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Azure ML Datasets makes it **easy to register your datasets** for use with your ML workloads

There will be various metadata associated to your dataset

You can upload new datasets and they will be **versioned**



Sample: Diabetes Version 1 (latest)

Details Consume Explore Models

Refresh Generate profile Unregister New version

Attributes Properties Tabular Description

Created by Andrew Brown

Web URL https://azureopendatastorage.blob.core.windows.net/mlsamples/diabetes.parquet

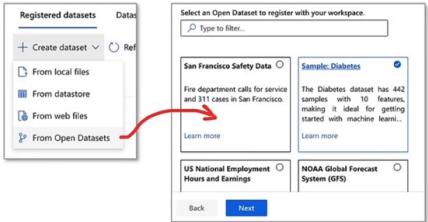
Profile Profile generation is running

Files in dataset 1 Total size of files in dataset 13.27 KB

# Azure Machine Learning Studio – Open Datasets

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

**Open Datasets** are publicly hosted datasets that are commonly used for learning how to build ML models



Azure has a curated list of open-datasets that you can quickly add to your data store. Great for learning how to use AutoML or Azure Machine Learning Designer

# Azure Machine Learning Studio – Pipelines

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

**Azure ML Pipelines** is an **executable workflow** of a complete machine learning task

Subtasks are encapsulated as a series of steps within the pipeline. Independent steps allow multiple data scientists to work on the same pipeline at the same time without over-taxing compute resources

Separate steps also make it easy to use different compute types/sizes for each step.

When you rerun a pipeline, the run jumps to the steps that need to be re-run, such as an updated training script.

Steps that do not need to be re-run are skipped

After a pipeline has been published, you can configure a REST endpoint, which allows you to rerun the pipeline from any platform or stack

You can build pipelines two ways:

- Using the Azure Machine Learning Designer
- **Programmatically using the Azure Machine Learning Python SDK**

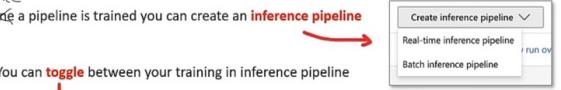
```
# In workspace, from conda  
blob_store = Datastore(ws, 'workspacedatastore')  
compute_target = ws.compute_targets['STANDARD_HD']  
experiment = Experiment(ws, "MyExperiment")  
  
Input_data = dataset.File.from_files(  
    blob_store, path='datasets/20newsgroups/20news.pckl')  
prepped_data = OutputFileDatasetConfig(name='output_path')  
  
data_prep_step = PythonScriptStep(  
    name='prep_data',  
    script_name='dataprep.py',  
    source_directory='scripts',  
    compute_target=compute_target,  
    arguments=[{"prepped_data_path": prepped_data.path},  
              {"input_dataset": Input_data}],  
    inputs=[Input_data],  
    outputs=[prepped_data])  
  
train_step = PythonScriptStep(  
    name='train',  
    script_name='train.py',  
    source_directory='scripts',  
    compute_target=compute_target,  
    arguments=[{"prepped_data_path": prepped_data.path},  
              {"source_directory": "train_src"}]  
    )  
steps = [data_prep_step, train_step]  
  
pipeline = Pipeline(workspace=ws, steps=steps)  
pipeline.run(experiment=experiment)  
pipeline_run.wait_for_completion()
```

# Azure Machine Learning Studio – Machine Learning Designer

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

Once a pipeline is trained you can create an **inference pipeline**

You can **toggle** between your training in inference pipeline



Training pipeline Real-time inference pipeline

Binary Classification with Feature Selection - Income Prediction-real time inference

Autosave on

Search by name, tags and description

99 assets in total

Datasets (1) Sample datasets (14) Data Input and Output (3) Data Transformation (15) Feature Selection (2) Statistical Functions (1) Machine Learning Algorithms (19) Model Training (4) Model Scoring & Evaluation (6) Python Language (2) R Language (1) Text Analytics (7) Computer Vision (6) Recommendation (5) Anomaly Detection (2)

Binary Classification with Feature Selection - Income Prediction-real time inference

Autosave on

Search by name, tags and description

99 assets in total

Datasets (1) Sample datasets (14) Data Input and Output (3) Data Transformation (15) Feature Selection (2) Statistical Functions (1) Machine Learning Algorithms (19) Model Training (4) Model Scoring & Evaluation (6) Python Language (2) R Language (1) Text Analytics (7) Computer Vision (6) Recommendation (5) Anomaly Detection (2)

Binary Classification with Feature Selection - Income Prediction-real time inference

Autosave on

Search by name, tags and description

99 assets in total

Datasets (1) Sample datasets (14) Data Input and Output (3) Data Transformation (15) Feature Selection (2) Statistical Functions (1) Machine Learning Algorithms (19) Model Training (4) Model Scoring & Evaluation (6) Python Language (2) R Language (1) Text Analytics (7) Computer Vision (6) Recommendation (5) Anomaly Detection (2)

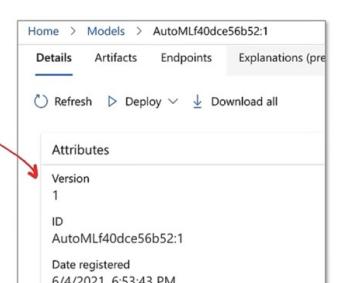
# Azure Machine Learning Studio – Models

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

**Model Registry** allows you to **create, manage and track your registered models** as incremental versions under the same name

Each time you register a model with the same name as an existing one, the registry assures that **it's a new version**.

Additionally, you can provide metadata tags and use the tags when you search for models.



Details	Artifacts	Endpoints	Explanations (pre)
Refresh Generate profile Unregister New version			
Attributes			
Version 1			
ID AutoML40dce56b52:1			
Date registered 6/4/2021, 6:53:43 PM			

# Azure Machine Learning Studio – Endpoints

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Azure ML Endpoints** allow you to **deploy machine learning models as a web service**

The workflow for deploying a model:

- Register the model
- Prepare an entry script
- Prepare an inference configuration
- Deploy the model locally to ensure everything works
- Choose a compute target
- Re-deploy the model to the cloud
- Test the resulting web service

## Realtime endpoints

An endpoint that provides remote access to invoke the ML model service running on either:

- Azure Kubernetes Service (AKS)
- Azure Container Instance (ACI)

## Pipeline endpoints

An endpoint that provide remote access to invoke an ML pipeline.

You can parametrize the pipeline endpoint for managed repeatability in batch scoring and retraining scenarios.

# Azure Machine Learning Studio – Endpoints

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

When you **deploy a model** to an endpoint it will either be deployed to:

- Azure Kubernetes Service (AKS)
- Azure Container Instance (ACI)

When you have deployed a real-time endpoint you can test the endpoint by sending a **single request** or a **batch request**.

The computing resource will not show in Azure Machine Learning Studio

You need to check AKS or ACI

# Azure Machine Learning Studio – Notebooks

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Azure has a built in **Jupyter-like Notebook editor** so you can build and train your ML models

```
By using Azure Machine Learning Compute, a managed service, data scientists can train machine learning models in clusters of virtual machines. Examples include VMs with GPU support. In this tutorial, you create Azure Machine Learning Compute as your training endpoint. You will submit Python code to run on this VM later in the tutorial. The code below creates the compute cluster for you if they don't already exist. Creation of compute takes approximately 5 minutes. If the Compute with that name is already in your workspace the code will skip the creation process.

# Import required libraries
from azureml.core import ComputeTarget, AmlCompute
from azureml.core.compute import ComputeTargetException

# choose a name for your cluster
cluster_name = "myamlcompute"
compute = ComputeTarget.get_by_name(ws, cluster_name)
if compute is None:
    # create the cluster
    compute = ComputeTarget.create(ws, cluster_name, AmlCompute.provisioning_config)

# Once the cluster is created, use get to inspect
# its status
compute.wait_for_completion(show_output=True)
```

## Choose Compute

You need to create a compute instance to run your Notebook

## Choose Kernel

You need to choose a Kernel which preload a programming language and programming libraries for different use cases

# Azure Machine Learning Studio – Notebooks

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

You can open the Notebook in a more familiar IDE:

- VSCode
- Jupyter Notebook (classic)
- **Jupyter Labs**

# Auto ML

## AutoML

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Automated machine learning (AutoML)** automates the process of creating an ML model.

With Azure AutoML you

- supply a dataset
- **Choose a Task Type** (Classification, Regression or Time Series Forecasting)
- Then AutoML will train and tune your model

### Classification

When you need to make a prediction based on several classes:

- binary classification: Yes or No
- multi-class classification: Red, Green, Blue

### Regression

When you need to predict a continuous number value

### Time Series Forecasting

When you need to predict the value based on time

## AutoML – Regression

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Regression** is a type of **supervised learning** in which **models learn using training data**, and apply those learnings to new data.

The goal of regression is to predict a variable in the future

## AutoML – Classification

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Classification** is a type of **supervised learning** in which **models learn using training data**, and apply those learnings to new data.

If you enable Deep Learning than you will likely want a **GPU compute**

The goal of classification models is to **predict which categories new data will fall into** based on learnings from its training data:

- **binary classification:** a record is labeled out of two possible labels eg: true or false
- **multiclass classification:** a record is labeled out of range of labels: happy, sad, mad, rad

## AutoML – Time Series Forecasting

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Forecast revenue, inventory, sales, or customer demand

An automated time-series experiment is treated as a **multivariate regression problem**

Past time-series values are "pivoted" to become additional dimensions for the regressor together with other predictor

unlike classical time series methods, has an advantage of naturally incorporating multiple contextual variables and their relationship to one another during training

### Time series forecasting

To predict values based on time

The time series forecasting method requires some additional information.

Time column \*

Select a time column...

Time series identifier(s)

Select column(s)...

Frequency \*

Autodetect

Forecast horizon \*

Autodetect

Enable deep learning

## AutoML – Data Guard Rails

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Data guardrails are run by Azure AutoML when **automatic featurization** is enabled.

A **sequence of checks** to **ensure high quality input data** is being used to train model.

## AutoML – Time Series Forecasting

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Advanced forecasting configuration includes:

- holiday detection and featurization
- time-series and DNN learners (Auto-ARIMA, Prophet, ForecastTCN)
- many models support through grouping
- rolling-origin cross validation
- configurable lags
- rolling window aggregate features

# AutoML – Automatic Featurization

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

During model training with AutoML, one of the following **scaling or normalization techniques** will be applied to each model

- StandardScaleWrapper — Standardize features by removing the mean and scaling to unit variance
- MinMaxScalar — Transforms features by scaling each feature by that column's minimum and maximum
- MaxAbsScaler — Scale each feature by its maximum absolute value
- RobustScalar — Scales features by their quantile range
- Principal component analysis (PCA) — Linear dimensionality reduction using Singular Value Decomposition of the data to project it to a lower dimensional space
- TruncatedSVDWrapper — This transformer performs linear dimensionality reduction by means of truncated singular value decomposition (SVD). Contrary to PCA, this estimator does not center the data before computing the singular value decomposition, which means it can work with scipy.sparse matrices efficiently
- SparseNormalizer — Each sample (that is, each row of the data matrix) with at least one non-zero component is rescaled independently of other samples so that its norm (l1 or l2) equals one

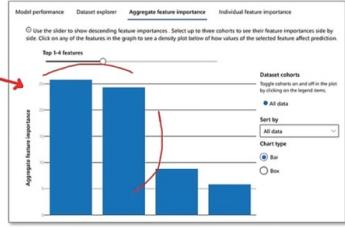
# AutoML – Explanation

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

ML Explainability (MLX) is the process of **explaining and interpreting** ML and deep learning models. MLX can help machine learning developers to better understand and interpret the model's behavior

After your top candidate model is selected by Azure AutoML you can get an explanation of the internals on various factors:

- Model Performance
- Dataset explorer
- Aggregate feature importance**
- Individual feature importance



# AutoML – Primary Metrics – Classification

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

## Classification Scenarios

- Suited for larger datasets that well-balanced
  - accuracy** — Image classification, Sentiment analysis, Churn prediction
  - average\_precision\_score\_weighted** — Sentiment analysis
  - norm\_macro\_recall** — Churn prediction
  - precision\_score\_weighted**
- Suited for small dataset that are imbalanced
  - AUC\_weighted** — Fraud detection, Image classification, Anomaly detection/spam detection

# AutoML – Model Selection

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

**Model selection** is the task of selecting a statistical model from a set of **candidate models**. Azure AutoML will use **many different ML Algorithms** and will recommend the best **performing candidate**

# AutoML – Primary Metrics

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

The primary metric parameter determines the metric to be used during model training for optimization.

## Classification

- accuracy
- AUC\_weighted
- average\_precision\_score\_weighted
- norm\_macro\_recall
- precision\_score\_weighted

## Regression and Time Series Forecasting

- spearman\_correlation
- normalized\_root\_mean\_squared\_error
- r2\_score
- normalized\_mean\_absolute\_error

# AutoML – Primary Metrics – Classification

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

## Regressions Scenarios

- Works well when value to predict encompasses a large range eg. 10K to 200K
  - spearman\_correlation**
  - r2\_score** — Airline delay, Salary estimation, Bug resolution time
- Works well when value to predict encompasses as smaller range eg. 10-20K
  - normalized\_root\_mean\_squared\_error** — Price prediction (house/product/tip), Review score prediction
  - normalized\_mean\_absolute\_error**

# AutoML – Primary Metrics – Time Series

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

## Time Series Scenarios

- Works well when value to predict encompasses a large range eg. 10K to 200K
  - spearman\_correlation**
  - r2\_score** — Price prediction (forecasting), Inventory optimization, Demand forecasting
- Works well when value to predict encompasses as smaller range eg. 10-20K
  - normalized\_root\_mean\_squared\_error** — Price prediction (forecasting), Inventory optimization, Demand forecasting
  - normalized\_mean\_absolute\_error**

# AutoML – Validation Type

Cheat sheets, Practice Exams and Flash cards [www.exampopro.co/ai-900](http://www.exampopro.co/ai-900)

**Model Validation** is when we **compare the results of our training dataset to our test dataset**.

Model Validation occurs after we train the model

With AutoML you can change the validation type

# Custom Vision



Custom Vision

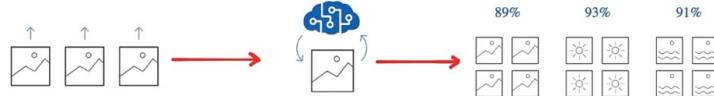
Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

Custom Vision is a **fully-managed no-code** service to quickly build your own **Classification and Object Detection ML models**.

This service is hosted on its own isolate domain at [www.customvision.ai](http://www.customvision.ai)

## Upload Images

Bring your own labeled images, or use Custom Vision to quickly add tags to any unlabeled images.



## Train

Use your labeled images to teach Custom Vision the concepts you care about.

## Evaluate

Use simple REST API calls to quickly tag images with your new custom computer vision model.

Within Custom Vision you setup projects and you need to select a **Project Type**

**Project Types** ⓘ

- Classification
- Object Detection

**Classification Types** ⓘ

- Multilabel (Multiple tags per image)
- Multiclass (Single tag per image)

## Classification

- Multi-label**
  - When we want to apply many tags to an image
  - Image contains both a Cat and a Dog
- Multi-class**
  - when we only have one possible tag to apply to an image:
  - It is either a Apple, Banana, Orange

## Object Detection

- When we detect various objects in an image

You will need to also choose a **Domain**

A Domain is a Microsoft Managed dataset that is used for training the ML model  
There are different domains that suited for different use cases



## Custom Vision – Image Classification Domains

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

### Classification Domains

**General** Optimized for a broad range of image classification tasks. If none of the other specific domains are appropriate, or if you're unsure of which domain to choose, select one of the General domains.

**General [A1]** Optimized for better accuracy with comparable inference time as General domain. Recommended for larger datasets or more difficult user scenarios. This domain requires more training time.

**General [A2]** Optimized for better accuracy with faster inference time than General[A1] and General domains. Recommended for most datasets. This domain requires less training time than General and General [A1] domains.

**Food** Optimized for photographs of dishes as you would see them on a restaurant menu. If you want to classify photographs of individual fruits or vegetables, use the Food domain.

**Landmark** Optimized for recognizable landmarks, both natural and artificial. This domain works best when the landmark is clearly visible in the photograph. This domain works even if the landmark is slightly obstructed by people in front of it.

**Retail** Optimized for images that are found in a shopping catalog or shopping website. If you want high-precision classifying between dresses, pants, and shirts, use this domain.

**Compact domains** Optimized for the constraints of real-time classification on edge devices.

### Object Detection Domains

#### General

Optimized for a broad range of object detection tasks. If none of the other domains are appropriate, or you are unsure of which domain to choose, select the General domain.

#### General [A1]

Optimized for better accuracy with comparable inference time as General domain. Recommended for more accurate region location needs, larger datasets, or more difficult user scenarios. This domain requires more training time, and results are not deterministic: expect a +1% mean Average Precision (mAP) difference with the same training data provided.

#### Logo

Optimized for finding brand logos in images.

#### Products on shelves

Optimized for detecting and classifying products on shelves.



## Custom Vision – Image Classification

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

For Image Classification you upload multiple images and you apply a single or multiple labels to the entire image.

The screenshot shows a workspace interface where multiple images are being tagged. A red arrow highlights the 'Tagged' tab under the 'Tags' section.

## Custom Vision – Training

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

When you are ready to train your model you have two options:

- Quick Training** – trains quickly but can be less accurate
- Advanced Training** – increase the compute time to improve your results

**Training Types** ⓘ

- Quick Training
- Advanced Training

**Advanced Training**

In most cases, the more time you select the better the model will be. You're charged based on the compute time used to train your model, so choose your budget based on your need.

Training budget: 1 hour | 1 hour | 24 hours

Send me an email notification after training completes

## Object Detection

## Classification

## Iterations

## Probability Threshold: 50%

## Overlap Threshold: 30%

## Iteration 1

## Training...

With each iteration of training our ML model will improve the evaluation metrics (**precision** and **recall**). The **probability threshold value** determines when to stop training when our evaluation metrics meet our desired threshold.

Once the **Classification** training job is complete we will get a report of the evaluation metrics **outcome**

#### Precision

- being exact and accurate
- select items that are relevant

#### Recall (Sensitivity or True Positive Rate)

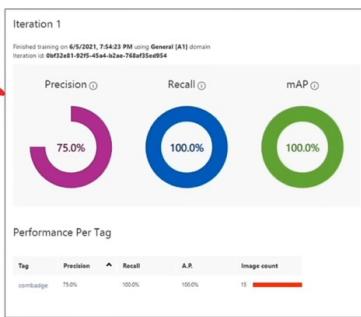
- How many relevant items returned

#### Average Precision (AP)



Once the **Object Detection** training job is complete we will get a report of the evaluation metrics **outcome**

- Precision
- Recall
- Mean Average Precision (mAP)



Before deploying our final trained model that can be invoked via an API Endpoint  
It is good practice to test our ML model using the **Quick Test** feature

To deploy our ML model to be accessible using our API Key and Endpoint we need to trigger the **Publish** action.

Once published we can get the **Prediction URL**

When you have a very large dataset you can use Smart Labeler to predict labels.  
Smart Labeler only works if you have trained the label

Smart labeler is when you want to increase your training set, and want to ML-assisted labeling to speed up this process.

## AI vs Gen AI

### Artificial Intelligence (AI)



AI refers to the development of computer systems that can **perform tasks typically requiring human intelligence**. These include **problem-solving, decision-making, understanding natural language, recognizing speech and images**, and more.



The primary goal of traditional AI is to create systems that can **interpret, analyze, and respond to human actions** or environmental changes efficiently and accurately. It aims to replicate or simulate human intelligence in machines.



AI applications are vast and include areas like **expert systems, natural language processing, speech recognition, and robotics**.



AI is used in various industries for tasks such as **customer service chatbots, recommendation systems in e-commerce, autonomous vehicles, and medical diagnosis**.

### Generative AI



Generative AI is a subset of AI that focuses on **creating new content or data** that is novel and realistic. It does not just interpret or analyze data but **generates new data itself**. It includes **generating text, images, music, speech, and other forms of media**.



It often involves advanced machine learning techniques, particularly deep learning models like **Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), and Transformer models (like GPT)**.



Generative AI is used in a range of applications including creating realistic **images and videos**, generating **human-like text**, composing **music**, creating virtual environments, and even drug discovery.



Examples: Tools like **GPT (Generative Pre-trained Transformer)** for text generation, **DALL-E** for image creation, and various deep learning models that compose music.

Feature	Artificial Intelligence (AI)	Generative AI
Functionality	Regular AI focuses on understanding and decision-making	Generative AI is about creating new, original outputs.
Data Handling	AI typically analyzes and makes decisions based on existing data	Generative AI uses existing data to generate new, unseen outputs.
Applications	Its applications span across various sectors, including data analysis, automation, natural language processing, and healthcare.	Its applications are more creative and innovative, focusing on content creation, synthetic data generation, deepfakes, and design.

# LLMs



## What is a Large Language Model (LLM)?

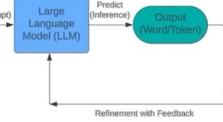
Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

A Large Language Model (LLM) such as GPT (Generative Pre-trained Transformer) works in a way that's similar to a complex, **automatic system that recognizes patterns and makes predictions**.

**Training on Large Datasets:** Initially, the model is trained on massive amounts of text data. This data can include **books, articles, websites, and other written material**.

During this training phase, the model learns patterns in language, such as grammar, word usage, sentence structure, and even style and tone.

**Understanding Context:** The model's design allows it to consider a wide context. This means it doesn't just focus on single words, but understands them in **relation to the words and sentences** that come **before and after**. This context understanding is important for generating coherent and relevant text.



## What is a Large Language Model (LLM)?

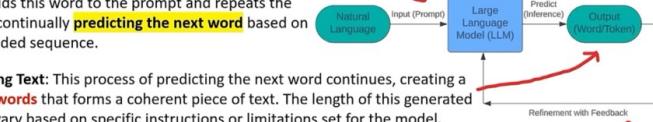
Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Predicting the Next Word:** When you give the model a **prompt** (a starting piece of text), it uses what it has learned to predict the next most likely word.

It then adds this word to the prompt and repeats the process, continually **predicting the next word** based on the extended sequence.

**Generating Text:** This process of predicting the next word continues, creating a **chain of words** that forms a coherent piece of text. The length of this generated text can vary based on specific instructions or limitations set for the model.

**Refinement with Feedback:** The model can be further **refined** and **improved** over time with **feedback**. This means it gets better at understanding and generating text as it is exposed to more data and usage.



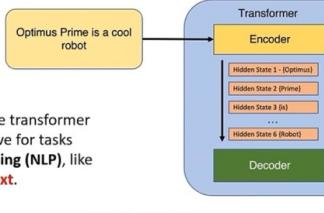
# Transformers



## Transformer models

A transformer model is a type of machine learning model that's especially good at **understanding and generating language**.

It's built using a structure called the transformer architecture, which is really effective for tasks involving **natural language processing (NLP)**, like **translating languages or writing text**.



Transformer model architecture consists of **two components**, or **blocks**:

1. **Encoder:** This part **reads and understands the input text**. It's like a smart system that goes through everything it's been taught (which is a lot of text) and picks up on the meanings of words and how they're used in different contexts.
2. **Decoder:** Based on what the encoder has learned, this part **generates new pieces of text**. It's like a skilled writer that can make up sentences that flow well and make sense.



## Tokenization

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Tokenization** in a transformer model is like turning a sentence into a puzzle. For example, you have the sentence: "I heard a dog bark loudly at a cat." To help a computer understand it, we chop up the sentence into pieces called **'tokens'**. Each piece can be a word or even a part of a word.

So, for our sentence, we give each word a number, like this:

- "I" might be 1
- "heard" might be 2
- "a" might be 3
- "dog" might be 4
- "bark" might be 5
- "loudly" might be 6
- "at" might be 7
- "a" is already tokenized as 3
- "cat" might be 8



Now, our sentence becomes a series of numbers: [1, 2, 3, 4, 5, 6, 7, 3, 8]. This is like giving each word a **special code**.

The computer uses these codes to **learn about the words and how they fit together**.

If a word repeats, like "a", we use its code again instead of making a new one.

As the computer reads more text, it keeps turning new words into new tokens with new numbers.

If it learns the word "meow," it might call it 9, and "skateboard" could be 10.



## Transformer models

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

There are different types of transformer models with specific jobs. For example:



BERT



GPT

**BERT** is good at **understanding the language**. It's like a librarian who knows where every book is and what's inside them. **Google** uses it to help its search engine understand what you're looking for.

**GPT** is good at **creating text**. It's like a skilled author who can write stories, articles, or conversations based on what it has learned.



## Embeddings

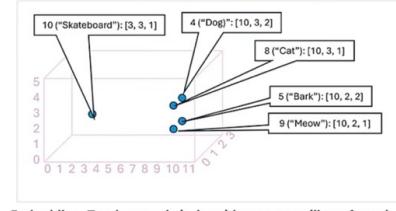
Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

To help a computer understand language, we turn words into tokens and then give each token a special **numeric code**, called an **embedding**. These embeddings are like a secret code that captures the meaning of the word. As a simple example, suppose the embeddings for our tokens consist of **vectors** with three elements, for example:

- 4 ("dog"): [10,3,2]
- 5 ("bark"): [10,2,2]
- 8 ("cat"): [10,3,1]
- 9 ("meow"): [10,2,1]
- 10 ("skateboard"): [3,3,1]

Words that have **similar meanings** or are used in similar ways get **codes that look alike**.

So, "dog" and "bark" might have similar codes because they are **related**.



This way, the computer can figure out which words are **similar to each other** just by looking at their codes. It's like giving each word a home on a map, and words that are neighbors on this map have related meanings.



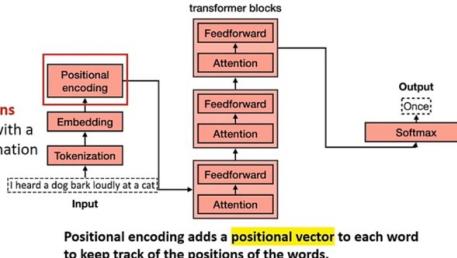
## Positional encoding

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Positional encoding** is a technique used to ensure that a language model, such as GPT (Generative Pre-trained Transformer) doesn't lose the **order of words** when processing natural language. This is important because the order in which words appear can change the meaning of a sentence.

Let's take the sentence "**I heard a dog bark loudly at a cat**" from our previous example:

Without positional encoding, if we simply tokenize this sentence and convert the **tokens** into **embedding vectors**, we might end up with a set of vectors that **lose the sequence information**.



## Positional encoding

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

However, by adding **positional encoding vectors to each word's embedding**, we ensure that each **position** in the sentence is uniquely identified:

- The embedding for "I" would be modified by adding a positional vector corresponding to position 1, labeled "I (1)".
- The embedding for "heard" would be altered by a vector for position 2, labeled "heard (2)".
- The embedding for "a" would be updated with a vector for position 3, labeled "a (3)", and reused with the same positional vector for its second occurrence.
- This process continues for each word/token in the sentence, with "dog (4)", "bark (5)", "loudly (6)", "at (7)", and "cat (8)" all receiving their unique positional encodings.

As a result, the sentence "**I heard a dog bark loudly at a cat**" is represented not just by a sequence of vectors for its words, but by a sequence of vectors that are influenced by the **position** of each word in the sentence.

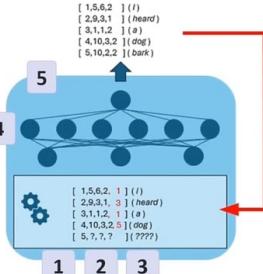
This means that even if another sentence had the same words in a different order, its overall representation would be different because the positional encodings would differ, reflecting the different sequence of words.

Attention in AI, especially in transformer models, is a way the model figures out how **important each word (or token) is to the meaning of a sentence**, particularly in **relation** to the other words around it. Let's reuse the sentence "I heard a dog bark loudly at a cat" to explain this better:

**Self-Attention:** Imagine each word in the sentence shining a flashlight on the other words. The brightness of the light shows how much one word should pay attention to the others when understanding the sentence. For "bark", the light might shine brightest on "dog" because they're closely related.

**Encoder's Role:** In the encoder part of a transformer model, attention helps decide **how to represent each word as a number (or vector)**. It's not just the word itself, but also its context that matters. For example, "bark" in "the bark of a tree" would have a different representation than "bark" in "I heard a dog bark", because the surrounding words are different.

- 1 **Token Embeddings:** Each word in the sentence is represented as a **vector** of numbers (its **embedding**).
- 2 **Predicting the Next Token:** The goal is to figure out what the next **word (token)** should be, also represented as a vector.
- 3 **Assigning Weights:** The attention layer looks at the sentence so far and decides how much **influence (weight)** each word should have on the next one.
- 4 **Calculating Attention Scores:** Using these weights, a new vector for the next token is calculated, which includes an attention score. **Multi-head attention** does this several times, focusing on different aspects of the words.
- 5 **Choosing the Most Likely Word:** A neural network takes these vectors with **attention scores** and picks the **word** from the vocabulary that most likely comes next.
- 6 **Adding to the Sequence:** The chosen word is **added to the existing sequence**, and the process **repeats for each new word**.



## Azure OpenAI Service

Azure OpenAI Service is a cloud-based platform designed to **deploy and manage advanced language models from OpenAI**. This service combines OpenAI's latest language model developments with the robust security and scalability of Azure's cloud infrastructure.

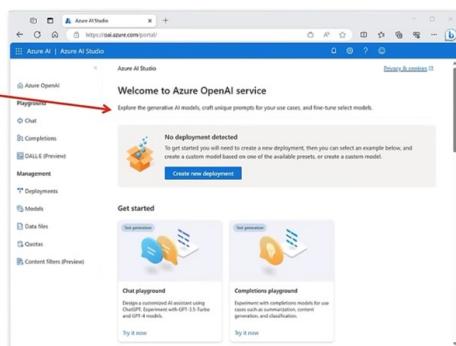
Azure OpenAI offers several **types of models** for different purposes:

- **GPT-4 Models:** These are the newest in the line of GPT models and can create **text and programming code** when given a prompt written in natural language.
- **GPT-3.5 Models:** Similar to GPT-4, these models also create text and code from natural language prompts. The GPT-3.5-turbo version is specially designed for **conversations**, making it a great choice for **chat applications and other interactive AI tasks**.
- **Embedding Models:** These models turn written **text into number sequences**, which is helpful for analyzing and comparing different pieces of text to find out how **similar** they are.
- **DALL-E Models:** These models can make **images from descriptions** given in words. The DALL-E models are still being tested and aren't shown in the Azure OpenAI Studio, so you don't have to set them up for use manually.

Developers can work with these models in **Azure OpenAI Studio**, a **web-based environment** where AI professionals can **deploy, test, and manage LLMs** that support generative AI app development on Azure.

Access is currently **limited** due to the high demand, upcoming product improvements, and Microsoft's commitment to responsible AI.

Presently, collaborations are being prioritized for those who already have a **partnership with Microsoft**, are engaged in lower-risk use cases, and are dedicated to including necessary safeguards.



### Pricing for Language models

Models	Context	Prompt (Per 1,000 tokens)	Completion (Per 1,000 tokens)
GPT-3.5-Turbo	4K	\$0.0015	\$0.002
GPT-3.5-Turbo	16K	\$0.003	\$0.004
GPT-3.5-Turbo-1106	16K	N/A	N/A
GPT-4-Turbo	128K	N/A	N/A
GPT-4-Turbo-Vision	128K	N/A	N/A
GPT-4	8K	\$0.03	\$0.06
GPT-4	32K	\$0.06	\$0.12

**Decoder's Role:** When generating new text, like completing a sentence, the decoder uses attention to figure out which words it already has are **most important** for **deciding what comes next**. If our sentence is "I heard a dog," the model uses **attention** to know that "heard" and "dog" are key to adding the next word, which might be "bark."

**Multi-Head Attention:** It's like having multiple flashlights, each **highlighting different aspects of the words**. Maybe one flashlight looks at the **meaning** of the word, another looks at its **role** in the sentence (like subject or object), and so on. This helps the model get a richer understanding of the text.

**Building the Output:** The decoder builds the sentence one word at a time, using **attention** at each step. It looks at the sentence so far, decides what's important, and then **predicts** the next word. It's an ongoing process, with each **new word influencing the next**.

Attention in transformer models is like a guide that helps the AI understand and create language by focusing on the most relevant parts of the text, considering both individual word meanings and their relationships within the sentence.

A transformer model like **GPT-4** works by taking a text **input (prompt)** and producing a well-structured **output (completion)**. During training, it learns from a vast array of text data, understanding how words are typically arranged in sentences.

The model knows the correct sequence of words but **hides (masks)** future words to learn how to **predict** them. When it tries to predict a word, it compares its guess to the actual word, gradually adjusting to reduce errors.

In practice, the model uses its training to assign **importance (weights)** to each word in a sequence, helping it guess the next word accurately. The result is that GPT-4 can create **sentences** that sound like they were written by a human.

However, this doesn't mean the model "knows" things or is "intelligent" in the human sense. It's simply very good at using its **large vocabulary and training to generate realistic text based on word relationships**.

Key concepts in using Azure OpenAI include **prompts and completions, tokens, resources, deployments, prompt engineering, and various models**:

**Prompts & Completions:** Users interact with the API by providing a text command in English, known as a **prompt**, and the model generates a text response, or completion.

• E.g., a prompt to count to five in a loop results in the model returning appropriate code.

**Tokens:** Azure OpenAI breaks down text into **tokens**, which are words or character chunks, to process requests. The number of tokens affects response latency and throughput.

• For images, token cost varies with image size and detail setting, with low-detail images costing fewer tokens and high-detail images costing more.

**Resources:** Azure OpenAI operates like other **Azure products** where users create a resource within their Azure Subscription.

**Deployments:** To use the service, users must deploy a **model via Deployment APIs**, choosing the specific model for their needs.

**Prompt Engineering:** Crafting **prompts** is crucial as they guide the model's output.

This requires skill, as prompt construction is nuanced and impacts the model's response.

**Models:** Various models offer different capabilities and pricing. **DALL-E** creates images from text, while **Whisper** transcribes and translates speech to text. Each has unique features suitable for different tasks.

In Azure AI Studio, you can deploy **large language models**, provide few-shot examples, and test them in Azure OpenAI Studio's Chat playground.

The image shows **Azure OpenAI's Chat playground interface**, where users can test and configure an AI chatbot.

In the middle, there's a **chat area** to type user messages and see the assistant's replies.

On the left, there's a menu for **navigation** and a section to set up the **assistant**, including a reminder to save changes.

On the right, adjustable **parameters** control the AI's response behavior, like **length, randomness, and repetition**. Users input queries, adjust settings, and observe how the AI responds to fine-tune its performance.

### Base models

Models	Usage per 1,000 tokens
Babbage-002	\$0.0004
Davinci-002	\$0.002

### Fine-tuning models

Models	Training per compute hour	Hosting per hour	Input Usage per 1,000 tokens	Output Usage per 1,000 tokens
Babbage-002	\$34	\$1.70	\$0.0004	\$0.0004
Davinci-002	\$68	\$3	\$0.002	\$0.002
GPT-3.5-Turbo	\$102	\$7	\$0.0015	\$0.002

### Image models

Models	Quality	Resolution	Price (per 100 Images)
Dall-E-3	Standard	1024 * 1024	\$4
	Standard	1024 * 1792, 1792 * 1024	\$8
Dall-E-3	HD	1024 * 1024	\$8
	HD	1024 * 1792, 1792 * 1024	N/A
Dall-E-2	Standard	1024 * 1024	\$2

### Embedding models

Models	Per 1,000 tokens
Ada	\$0.0001

### Speech models

Models	Per hour
Whisper	\$0.36

# Copilots

## What are copilots?

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](https://www.exampro.co/ai-900)

**Copilots** are a new type of computing tool that integrates with applications to help users with **common tasks using generative AI models**. They are designed using a standard architecture, allowing developers to create **custom copilots** tailored to specific business needs and applications.

- Copilots might appear as a chat feature beside your document or file, and they **utilize the content within the product to generate specific results**.

Creating a copilot involves several steps:

1. Training a **large language model** with a vast amount of data.
2. Utilizing services like **Azure OpenAI Service**, which provide pre-trained models that developers can either use as-is or fine-tune with their own data for more specific tasks.
3. **Deploying** the model to make it available for use within applications.
4. **Building copilots** that **prompt** the models to generate usable content.
5. Business users can use copilots to boost their **productivity and creativity** with AI-generated content.

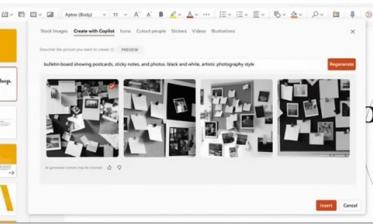
Copilots have the potential to revolutionize the way we work. These copilots use generative AI to help with first drafts, information synthesis, strategic planning, and much more.

## Copilot Examples

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](https://www.exampro.co/ai-900)

**Microsoft 365 Copilot** is designed to be a partner in your workflow, integrated with productivity and communication tools like **PowerPoint** and **Outlook**.

- It's there to help you craft effective **documents**, **design spreadsheets**, put together **presentations**, manage emails, and streamline other tasks.



## GitHub Copilot

**GitHub Copilot** is a tool that helps software developers, offering real-time assistance as they **write code**. It offers more than suggesting code snippets; it can help in **documenting the code** for better understanding and maintenance.

- Copilot also helps **test out code**, which means coders can work smarter and make fewer mistakes.

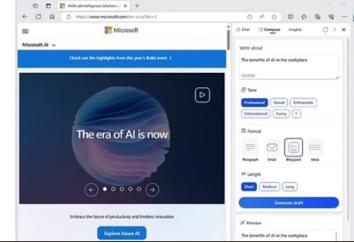
```
1 // Users > benyone > JS Testjs > findimagesWithoutAlt
2 // find all Images without alternate text
3 // and give them a border
4 function findimagesWithoutAlt() {
5     var Images = document.getElementsByTagName("img");
6     for (var i = 0; i < Images.length; i++) {
7         if (Images[i].alt == "") {
8             Images[i].style.border = "2px solid red";
9         }
10    }
11 }
```

## Copilot Examples

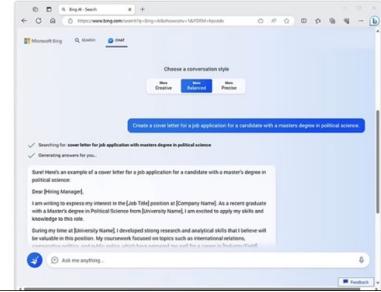
Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](https://www.exampro.co/ai-900)

**Microsoft Copilot** is integrated into various applications to assist users in creating **documents**, **spreadsheets**, **presentations**, and **more**, by generating content, summarizing information, and aiding in strategic planning.

- It is used across Microsoft's suite of products and services to enhance user experience and efficiency.



The Microsoft Bing search engine provides a copilot to help when **browsing** or **searching** the Internet by generating **natural language answers** to questions based on context rather than just search results of indexed pages.



# Prompt Engineering

## Prompt engineering

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](https://www.exampro.co/ai-900)

**Prompt engineering** is a process that **improves the interaction** between humans and generative AI. It involves **refining the prompts** or instructions given to an AI application to **generate higher quality responses**. This process is valuable for both the developers who create AI-driven applications and the end-users who interact with them.

For example, developers may build a generative AI application for teachers to create multiple-choice questions related to text students read. During the development of the application, developers can **add other rules** for what the program should do with the prompts it receives.

### System messages

Prompt engineering techniques include defining a system message. The message sets the **context for the model** by describing **expectations** and **constraints**.

For example, "You're a helpful assistant that responds in a **cheerful, friendly manner**".

These system messages determine constraints and styles for the model's responses.

### Writing good prompts

To maximize the utility of AI responses, it is essential to be **precise and explicit** in your prompts.

A well-structured prompt, such as "**Create a list of 10 things to do in Edinburgh during August**," directs the AI to produce a targeted and relevant output, achieving better results.

## Prompt engineering workflow

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](https://www.exampro.co/ai-900)

1. **Task Understanding:** Know what you want the AI to do.

2. **Crafting Prompts:** Write instructions for the AI.

3. **Prompt Alignment:** Make sure instructions match what the AI can do.

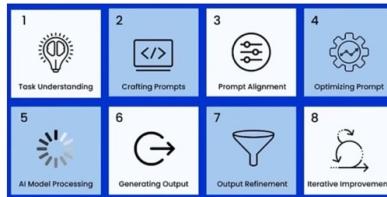
4. **Optimizing Prompt:** Improve the instructions for better AI responses.

5. **AI Model Processing:** The AI thinks about the instructions.

6. **Generating Output:** The AI gives an answer or result.

7. **Output Refinement:** Fix or tweak the AI's answer.

8. **Iterative Improvement:** Keep improving the instructions and answers.

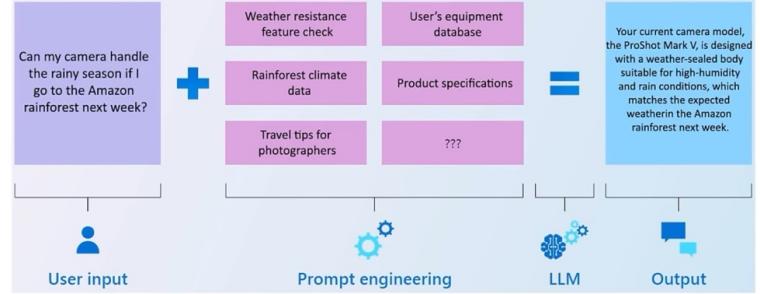


## Prompt engineering

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](https://www.exampro.co/ai-900)

**Zero-shot learning** refers to an AI model's ability to correctly perform a task **without any prior examples or training** on that specific task.

**One-shot learning** involves the AI model learning from a **single example or instance** to perform a task.



# Grounding

## Grounding

Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Grounding** in prompt engineering is a technique used in large language models (LLMs) where you provide **specific, relevant context within a prompt**. This helps the AI to produce a more accurate and related response.

For example, if you want an LLM to summarize an email, you would include the actual email text in the prompt along with a command to summarize it. This approach allows you to leverage the LLM for tasks it wasn't explicitly trained on, without the need for retraining the model.

### Prompt engineering vs Grounding

**Prompt engineering** broadly refers to the art of crafting effective prompts to produce the desired output from an AI model. **Grounding** specifically involves enriching prompts with relevant context to improve the model's understanding and responses.

**Grounding** ensures the AI has **enough information to process the prompt** correctly, whereas **prompt engineering** can also include techniques like format, style, and the strategic use of examples or questions to guide the AI.

## Grounding options

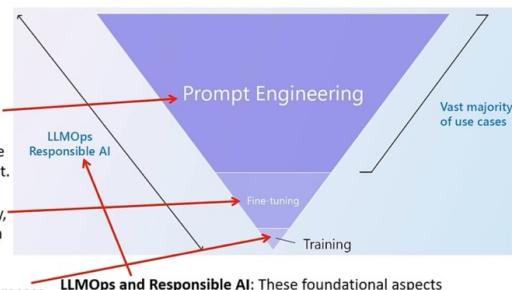
Cheat sheets, Practice Exams and Flash cards [www.exampro.co/ai-900](http://www.exampro.co/ai-900)

**Grounding Options:** These are techniques to ensure LLM outputs are accurate and adhere to responsible AI principles.

**Prompt Engineering:** Placed at the top, indicating its broad applicability, this involves designing prompts to direct the AI toward generating the desired output.

**Fine-Tuning:** A step below in complexity, where LLMs are trained on specific data to improve their task performance.

**Training:** The most resource-intensive process, at the triangle's base, suggesting its use in more extensive customization needs.



**LLMops and Responsible AI:** These foundational aspects emphasize the importance of operational efficiency and ethical standards across all stages of LLM application development.

Vast majority of use cases