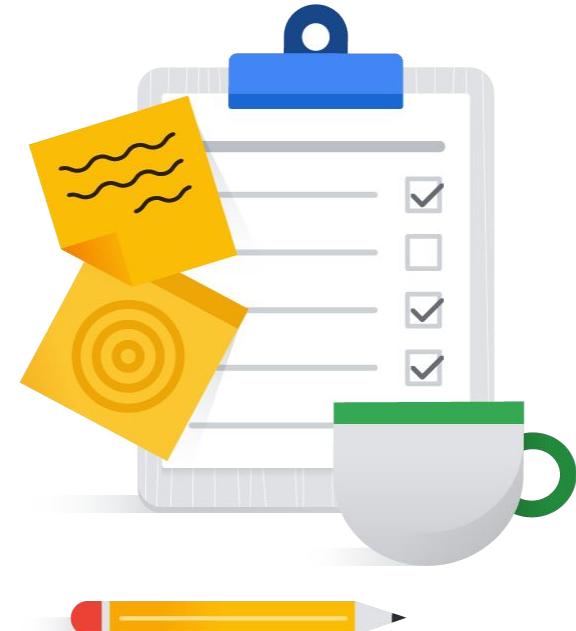




# Model Garden

# Model Garden

- 01 [Introduction to Model Garden](#)
- 02 Available Model Types
- 03 Working with Gen AI Studio
- 04 Working with Model Registry
- 05 Intro: Course Use Cases
- 06 Lab: Exploring Model Garden



# Model Garden

- 01 Introduction to Model Garden
- 02 Available Model Types
- 03 Working with Gen AI Studio
- 04 Working with Model Registry
- 05 Intro: Course Use Cases



**Modalities**

Language 39

Vision 79

Tabular 2

Documents 2

Speech 1

Video 4

**Tasks**

Generation 39

Classification 42

Detection 29

Extraction 8

Recognition 8

Translation 4

Embedding 3

Segmentation 6

Retrieval 1

Tracking 1

**Features**

Generative AI Studio 7

API available 0

Search models

Suggestions text embedding essay outline BERT

Pre-trained multi-task models that can be further tuned or customised for specific tasks. Models marked with ◆ are available in [Generative AI Studio](#).

◆ Generative AI Language

PaLM 2 for Text

Fine-tuned to follow natural language instructions and is suitable for a variety of language tasks, such as: classification, extraction, summarization and content...  
text-bison

[VIEW DETAILS](#)

Foundation Language

Llama 2

Fine-tune & deploy Meta's Llama 2 models on Vertex AI.

Llama2-7B-001

[VIEW DETAILS](#)

Foundation Language

Embeddings for Text

Text embedding is an important NLP technique that converts textual data into numerical vectors that can be processed by machine learning algorithms, especially large models...  
textembedding-gecko@001

[VIEW DETAILS](#)

◆ Generative AI Vision

Imagen for Image Generation

Create or edit studio-grade images at scale via text prompts

google/imagegeneration-002

[VIEW DETAILS](#)[▼ SHOW ALL \(54\)](#)**Fine-tunable models**

Models that data scientists can further fine-tune through a custom notebook or pipeline.

Classification Vision

tfhub/EfficientNetV2

EfficientNet V2 are a family of image classification models, which achieve better parameter efficiency and faster training speed than prior arts.

tensorflow-hub/efficientnetv2

[VIEW DETAILS](#)

Classification Vision

tfvision/vit

The Vision Transformer (ViT) is a transformer-based architecture for image classification.

tfvision/vit-s16

[VIEW DETAILS](#)

Detection Vision

tfvision/SpineNet

SpineNet is an image object detection model generated using Neural Architecture Search.

tfvision/spinenet49

[VIEW DETAILS](#)

Detection Vision

tfvision/YOLO

YOLO algorithm is a one-stage object detection algorithm that can achieve real-time performance on a single GPU.

tfvision/scaled-yolo

[VIEW DETAILS](#)

## Model Garden

[EXPLORE GENERATIVE AI](#)
 Search models

Suggestions

[text embedding](#)
[essay outline](#)
[BERT](#)

Pre-trained

## Modalities

Language 39

Vision 79

Tabular 2

Documents 2

Speech 1

Video 4

## Tasks

Generation 39

Classification 42

Detection 29

Extraction 8

Recognition 8

Translation 4

Embedding 3

Segmentation 6

Retrieval 1

Tracking 1

## Features

Generative AI Studio 7

API available 0

 Generative AI

Language

## PaLM 2 for Text

Fine-tuned to follow natural language instructions and is suitable for a variety of language tasks, such as: classification, extraction, summarization and content...  
text-bison

[VIEW DETAILS](#)
 Foundation

Language

## Llama 2

Fine-tune & deploy Meta's Llama 2 models on Vertex AI.

Llama2-7B-001

[VIEW DETAILS](#)
 Foundation

Language

## Embeddings for Text

Text embedding is an important NLP technique that converts textual data into numerical vectors that can be processed by machine learning algorithms, especially large models...  
textembedding-gecko@001

[VIEW DETAILS](#)
 Generative AI

Vision

## Imagen for Image Generation

Create or edit studio-grade images at scale via text prompts

google/imagegeneration-002

[VIEW DETAILS](#)[▼ SHOW ALL \(54\)](#)

## Fine-tunable models

Models that data scientists can further fine-tune through a custom notebook or pipeline.

 Classification

 Vision

## tfhub/EfficientNetV2

EfficientNet V2 are a family of image classification models, which achieve better parameter efficiency and faster training speed than prior arts.

tensorflow-hub/efficientnetv2

[VIEW DETAILS](#)
 Classification

 Vision

## tfvision/vit

The Vision Transformer (ViT) is a transformer-based architecture for image classification.

tfvision/vit-s16

[VIEW DETAILS](#)
 Detection

 Vision

## tfvision/SpineNet

SpineNet is an image object detection model generated using Neural Architecture Search.

tfvision/spinenet49

[VIEW DETAILS](#)
 Detection

 Vision

## tfvision/YOLO

YOLO algorithm is a one-stage object detection algorithm that can achieve real-time performance on a single GPU.

tfvision/scaled-yolo

[VIEW DETAILS](#)

Model Type	Count	Description
Language	39	Language models
Vision	79	Vision models
Tabular	2	Tabular models
Documents	2	Document models
Speech	1	Speech models
Video	4	Video or Text models
Tasks		Task models that can be further tuned or customized for specific tasks.
Generation	39	Follow natural language and is suitable for a range of language tasks, such as: extraction, translation and content...
Classification	42	Classification models
Detection	29	Detection models
Extraction	8	Extraction models
Recognition	8	Recognition models
Translation	4	Translation models
Embedding	3	Embedding models
Segmentation	6	Segmentation models
Retrieval	1	Retrieval models
Tracking	1	Tracking models
MY MODELS		
Text embedding		
essay outline		
BERT		
Sample models		
PaLM 2 for Chat	Foundation Language	Fine-tuned to conduct natural conversation. Use this model to build and customize your own chatbot application. chat-bison@001
Embeddings for text	Foundation Language	Text embedding is an important NLP technique that converts textual data into numerical vectors that can be processed by machine learning algorithms, especially large models... textembedding-gecko@001
Chirp	Foundation Speech	Chirp is a version of a Universal Speech Model that has over 2B parameters and can transcribe in over 100 languages in a single model. chirp-rnnt1
Embeddings for Image	Foundation Vision	Generates vectors based on images, which can be used for downstream tasks like image classification, image search, and so on. imageembedding-001
tfvision/vit	Classification Vision	The Vision Transformer (ViT) is a transformer-based architecture for image classification. tfvision/vit-s16
tfvision/SpineNet	Detection Vision	SpineNet is an image object detection model generated using Neural Architecture Search. tfvision/spinenet49
tfvision/YOLO	Detection Vision	YOLO algorithm is a one-stage object detection algorithm that can achieve real-time performance on a single GPU. tfvision/scaled-yolo
Stable Diffusion Inpainting	Foundation Language Vision	Stable Diffusion Inpainting is a latent diffusion model capable of inpainting images given any text input and a mask image. runwayml/stable-diffusion-inpainting

**Modalities**

Language 39

Vision 79

Tabular 2

Documents 2

Speech 1

Video 4

**Tasks**

Generation 39

Classification 42

Detection 29

Extraction 8

Recognition 8

Translation 4

Embedding 3

Segmentation 6

Retrieval 1

Tracking 1

**Features**

Generative AI Studio 7

API available 0

## Foundation models

Pre-trained multi-task models that can be further tuned or customised for specific tasks. Models marked with ◆ are available in [Generative AI Studio](#).

◆ Generative AI Language

### PaLM 2 for Text

Fine-tuned to follow natural language instructions and is suitable for a variety of language tasks, such as: classification, extraction, summarization and content...  
text-bison

[VIEW DETAILS](#)

Foundation Language

### Llama 2

Fine-tune & deploy Meta's Llama 2 models on Vertex AI.

Llama2-7B-001

[VIEW DETAILS](#)

Foundation Language

### Embeddings for Text

Text embedding is an important NLP technique that converts textual data into numerical vectors that can be processed by machine learning algorithms, especially large models...  
textembedding-gecko@001

[VIEW DETAILS](#)

#### Fine-tuned

#### Models that can be fine-tuned

#### Classified

▼ SHOW ALL (54)

tfhub/EfficientNetV2

EfficientNet V2 are a family of image classification models, which achieve better parameter efficiency and faster training speed than prior arts.

[tensorflow-hub/efficientnetv2](#)

[VIEW DETAILS](#)

tfvision/vit

The Vision Transformer (ViT) is a transformer-based architecture for image classification.

[tfvision/vit-s16](#)

[VIEW DETAILS](#)

tfvision/SpineNet

SpineNet is an image object detection model generated using Neural Architecture Search.

[tfvision/spinenet49](#)

[VIEW DETAILS](#)

tfvision/YOLO

YOLO algorithm is a one-stage object detection algorithm that can achieve real-time performance on a single GPU.

[tfvision/scaled-yolo](#)

[VIEW DETAILS](#)

# Model Card

## PaLM 2 for Text

Fine-tuned to follow natural language instructions and is suitable for a variety of language tasks, such as: classification, extraction, summarization and content generation.

[OPEN PROMPT DESIGN](#)
[VIEW API CODE](#)
[OVERVIEW](#)
[USE CASES](#)
[DOCUMENTATION](#)

### Overview

**text-bison** is the name of the PaLM 2 for text large language model that understands and generates language. It's a foundation model that performs well at a variety of natural language tasks such as sentiment analysis, entity extraction, and content creation. The type of content that **text-bison** can create includes document summaries, answers to questions, and labels that classify content.

The PaLM 2 for text is ideal for tasks that can be completed with one API response, without the need for continuous conversation. For text tasks that require back-and-forth interactions, use the [PaLM 2 for chat](#).

### Use cases

- Summarization:** Create a shorter version of a document that incorporates pertinent information from the original text. For example, you might want to summarize a chapter from a textbook. Or, you could create a succinct product description from a long paragraph that describes the product in detail.
- Question answering:** Provide answers to questions in text. For example, you might automate the creation of a Frequently Asked Questions (FAQ) document from knowledge base content.
- Classification:** Assign a label to provided text. For example, a label might be applied to text that describes how grammatically correct it is.

### Resource ID

text-bison@001

### Tags

#### Task

Classification  
Detection  
Extraction  
Generation  
Recognition  
Translation

### Versions

Resource ID	Release date	Release stage	Description
text-bison	2023-08-29	General Availability	
text-bison-32k	2023-08-29	Public Preview	
text-bison@001	2023-06-07	General Availability	
text-bison@001	2023-05-10	Public Preview	Quality improvements and restage -001 as the first stable base model release

# Model Garden

## Foundation

Google's state-of-the-art multimodal models

## Pre-trained APIs

Google's pre-trained, task-specific models

## Open Source

Enterprise-ready open source models

## Third-Party

Coming soon!

# Foundation Models

Across a variety of model sizes to address use cases



## PaLM for Text

Custom language tasks



## PaLM for Chat

Multi-turn conversations with session context



## Imagen for Text to Image

Create and edit images from simple prompts



## Embeddings API for Text and Image

Extract semantic information from unstructured data



## Chirp for Speech to Text

Build voice enabled applications



## Codey for Code Generation

Improve coding and debugging

# PaLM 2

## Text

Perform language tasks like summarization, question answering, and classification from natural language prompts.

## Chat

Text tasks that flow like natural, multi-turn conversations such as help agents.

## Embeddings

Retrieve the vector representation for the semantic meaning of input text.

# Codey

## Code Generation

Generate functions,  
tests, and object classes  
from natural language  
descriptions.

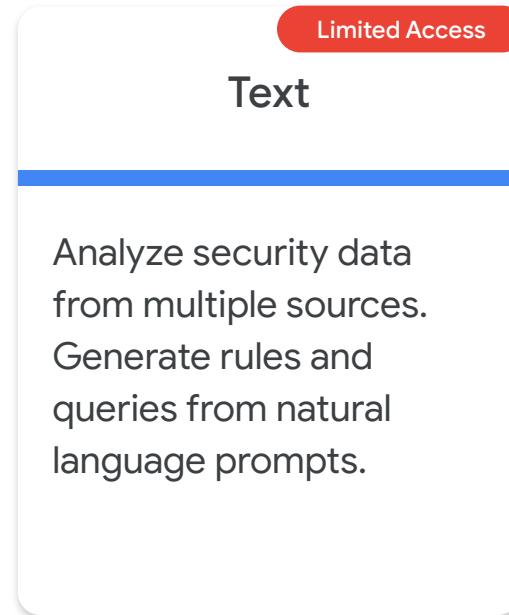
## Code Completion

In-line code generation  
based on surrounding  
code.

## Code Chat

Multi-turn reasoning,  
debugging, and  
explanations of code.

# Sec-PaLM2



# Imagen

Limited Access

## Image Generation

Create images from text prompts.

## Embeddings

Vector representation of an image. Also generates text vectors for querying text-to-image and image-to-text.

# Chirp

## Transcription

Convert speech to text  
for over 100 languages.

# Other foundation models

## Image Generation

- Stable Diffusion
- Stable Diffusion Inpainting
- Stable Diffusion LoRA
- InstructPix2Pix
- ControlNet
- FaceStylizer

## Image Processing

- PaLI zero-shot
- BLIP 2
- CLIP
- OWL-ViT
- ViT GPT2
- LayoutLM for VQA
- BiomedCLIP
- ImageBind

## Natural Language

- BERT
- T5-FLAN
- NLLB (Translation)
- Mistral-7B
- RoBERTa-large (PEFT)
- Palmyra Med

# Fine-tunable models

## Foundation

- PaLM (via Gen AI Studio)
- Stable Diffusion Inpainting
- ControlNet

## Image Classification

- EfficientNetV2
- ResNet
- BEiT
- DeiT
- MobileNet
- ... many more!

## Object Detection

- SpineNet
- YOLO
- Faster R-CNN
- Mask R-CNN
- RetinaNet

# Task-specific solutions: Vision

## Detection

- Occupancy
- Person/Vehicle
- PPE
- Objects
- Person blur
- Watermark
- Faces

## Recognition

- Product
- Tag

## Extraction

- Text
- Form Parser
- Document AI
- OCR

## Classification

- Moderation

# Task-specific solutions: Language

## Classification

- Entity analysis
- Sentiment analysis
- Entity sentiment analysis
- Text Moderation

## Translation

- Text Translation

## Extraction

- Syntax analysis

# Task-specific solutions: Tabular

## Classification / Regression

- AutoML E2E
- TabNet

Sometimes, you don't need a model.  
You just need a **prompt**.

# Model Garden

- 01 Introduction to Model Garden
- 02 Available Model Types
- 03 **Working with Gen AI Studio**
- 04 Working with Model Registry
- 05 Intro: Course Use Cases



# Generative AI Studio

## PaLM 2 for Text

Fine-tuned to follow natural language instructions and is suitable for a variety of language tasks, such as: classification, extraction, summarization and content generation.

[OPEN PROMPT DESIGN](#)[VIEW API CODE](#)[OVERVIEW](#)[USE CASES](#)[DOCUMENTATION](#)

### Overview

**text-bison** is the name of the PaLM 2 for text large language model that understands and generates language. It's a foundation model that performs well at a variety of natural language tasks such as sentiment analysis, entity extraction, and content creation. The type of content that **text-bison** can create includes document summaries, answers to questions, and labels that classify content.

The PaLM 2 for text is ideal for tasks that can be completed with one API response, without the need for continuous conversation. For text tasks that require back-and-forth interactions, use the [PaLM 2 for chat](#).

### Use cases

- **Summarization:** Create a shorter version of a document that incorporates pertinent information from the original text. For example, you might want to summarize a chapter from a textbook. Or, you could create a succinct product description from a long paragraph that describes the product in detail.
- **Question answering:** Provide answers to questions in text. For example, you might automate the creation

## Codey for Code Generation

Generates code based on natural language input. Good for writing functions, classes, unit tests, and more.

[OPEN PROMPT DESIGN](#)[VIEW API CODE](#)[OVERVIEW](#)[USE CASES](#)[DOCUMENTATION](#)

### Overview

**code-bison** is the name of the model that supports code generation. It's a foundation model that generates code based on a natural language description. The type of content that **code-bison** can create includes functions, web pages, and unit tests. **code-bison** is supported by the code generation Codey API. The Codey APIs are in the PaLM API family.

The Codey API for code generation is ideal for tasks that can be completed with one API response, without the need for continuous conversation. For code tasks that require back-and-forth interactions, use the [Codey API for code chat](#).

### Use cases

Some common used cases for code generation are:

- **Unit tests:** Use the prompt to request a unit test for a function.
- **Write a function:** Pass a problem to the model to get a function that solves that problem.
- **Create a class:** Use a prompt to describe the purpose of a class and have code that defines the class returned.

**Prompt**

Write a prompt and then click Submit

**Response**

Markdown

The model will generate a response after you click Submit

Model may display inaccurate or offensive information that doesn't represent Google's view. Not all languages are supported. [Learn more.](#)

Add stop sequence

Press Enter after each sequence

Streaming responses

Print responses as they're generated

Safety filter threshold

Block few

Submit

RESET PARAMETERS

Google Cloud

Region

us-central1 (Iowa)



Model

text-bison (latest)



Temperature

0

1

0.2

Advanced

Token limit

1

2048

1024

Top-k

1

40

40

Top-p

0

1

0.8

Max. responses

1

8

1

### Prompt

Write a prompt and then click Submit



### Response

Markdown

The model will generate a response after you click Submit

Model may display inaccurate or offensive information that doesn't represent Google's view. Google models currently only support English.

We want your feedback.

### Model

Filter Type to filter

### PaLM 2

- text-bison (latest)**  PREVIEW  
Latest
- text-bison-32k (latest)**  PREVIEW  
Latest
- text-bison@001** Best value

### Codey (PaLM 2)

- code-bison (latest)**  PREVIEW

### REPORT INAPPROPRIATE RESPONSES

Content processed through Vertex AI is assessed against a list of safety attributes. Confidence scores for these attributes are returned in API responses [but are not visible in the UI](#).

**Prompt**

Write a prompt and then click Submit

**Response**

Markdown

The model will generate a response after you click Submit

Model may display inaccurate or offensive information that doesn't represent Google's view. Google models currently only support English.

**Temperature**

0  1

0.2

**Advanced****Token limit**

0  1 2048

1024

**Top-k**

1  40

40

**Top-p**

0  1

0.8

**Max. responses**

1  8

1

**Add stop sequence**

Press Enter after each sequence

**Streaming responses**

Print responses as they're generated

**Safety filter threshold**

Block few

**Context** ?

Write your context here...

**Examples** ?

IMPORT

ADD INPUT COLUMN

**INPUT** Add an input prefix**OUTPUT** Add an output prefix

Write an input to add a new row

Write an output to add a new row

**Test** ?

Batch testing

ADD TO EXAMPLES

**INPUT** Add an input prefix**OUTPUT** Add an output prefix

Markdown

Write an input

Model responses will appear here

Model may display inaccurate or offensive information that doesn't represent Google's view. Not all languages are supported. [Learn more](#) Region ?

us-central1 (Iowa)

Model ?

text-bison (latest)

Temperature ?

0 1

0.2

**Advanced**Token limit ?

1 2048

1024

Top-k ?

1 40

40

Top-p ?

0 1

0.8

Max. responses ?

1 8

1

Add stop sequence ?

Press Enter after each sequence

Streaming responses

Print responses as they're generated

Safety filter threshold ?

Block few

Context ?

You are a machine learning practitioner training a junior staff member. Provide answers to questions that include practical use cases to help the new employee learn.

Examples ?[ADD INPUT COLUMN](#)Test ?

INPUT Add an input prefix

Why would I choose precision versus recall as my evaluation metric?

OUTPUT Add an output prefix

[ADD TO EXAMPLES](#) Markdown

Precision and recall are two important metrics for evaluating the performance of a binary classifier. Precision is the fraction of positive predictions that are actually positive, while recall is the fraction of actual positives that are predicted as positive.

Precision is important when the cost of false positives is high. For example, if you are using a classifier to identify fraudulent transactions, you would want to have a high precision so that you don't accidentally block legitimate transactions.

Recall is important when the cost of false negatives is high. For example, if you are using a classifier to identify cancer, you would want to have a high recall so that you don't miss any cases of cancer.

The choice of which metric to use depends on the specific application. In general, you should choose the metric that is most important for your particular application.

Token limit ?

1  2048

1024

Top-k ?

1  40

40

Top-p ?

0  1

0.8

Max. responses ?

1  8

1

Add stop sequence ?

Press Enter after each sequence

 Streaming responses

Print responses as they're generated

Safety filter threshold ?

Block few

[SUBMIT](#)[RESET PARAMETERS](#)

&gt; Cloud

## Context ?

You are a machine learning practitioner training a junior staff member. Provide answers to questions that include practical use cases to help the new employee learn.

## Examples ?

ADD INPUT COLUMN

## Test ?

ADD TO EXAMPLES

INPUT Add an input prefix

OUTPUT Add an output prefix

Markdown

Why would I choose precision versus recall as my evaluation metric?

Precision and recall are two important metrics for evaluating the performance of a binary classifier. Precision is the fraction of positive predictions that are actually positive, while recall is the fraction of actual positives that are predicted as positive.

Precision is important when the cost of false positives is high. For example, if you are using a classifier to identify fraudulent transactions, you would want to have a high precision so that you don't accidentally block legitimate transactions.

Recall is important when the cost of false negatives is high. For example, if you are using a classifier to identify cancer, you would want to have a high recall so that you don't miss any cases of cancer.

The choice of which metric to use depends on the specific application. In general, you should choose the metric that is most important for your particular application.

## Model

text-bison (latest)

## Temperature ?

0  1

0.2

## Advanced

## Token limit ?

1  2048

1024

## Top-k ?

1  40

40

## Top-p ?

0  1

0.8

## Max. responses ?

1  8

1

Add stop sequence

?

Press Enter after each sequence

 Streaming responses

Print responses as they're generated

## Safety filter threshold

Block few

SUBMIT

RESET PARAMETERS

**Context** ?

You are a machine learning practitioner training a junior staff member. Provide answers to questions that include practical use cases to help the new employee learn.

**Examples** ?**ADD INPUT COLUMN****Test** ?**ADD TO EXAMPLES****INPUT** Add an input prefix**OUTPUT**

Add an output prefix

**Markdown**

Why would I choose precision versus recall as my evaluation metric?

Precision and recall are two important metrics for evaluating the performance of a binary classifier. Precision is the fraction of positive predictions that are actually positive, while recall is the fraction of actual positives that are predicted as positive.

Precision is important when the cost of false positives is high. For example, if you are using a classifier to identify fraudulent transactions, you would want to have a high precision so that you don't accidentally block legitimate transactions.

Recall is important when the cost of false negatives is high. For example, if you are using a classifier to identify cancer, you would want to have a high recall so that you don't miss any cases of cancer.

The choice of which metric to use depends on the specific application. In general, you should choose the metric that is most important for your particular application.

**Model** -

text-bison (latest)

Temperature ?



0.2

**Advanced**

Token limit ?



1024

Top-k ?



40

Top-p ?



0.8

Max. responses ?



1

Add stop sequence ?

Press Enter after each sequence

Streaming responses

Print responses as they're generated

Safety filter threshold

Block few

**SUBMIT****RESET PARAMETERS**

# Tune a Model

[Create a tuned model](#)

- 1 Tuning type
- 2 Model details
- 3 Tuning dataset
- 4 Evaluation (optional)

START TUNING

## Choose a tuning method

Tuning improves model quality for a specific domain or data set. The recommended tuning method depends on the data that you have available, your goals and use case.

[Learn more about tuning](#)

Supervised tuning

Uses example prompt and model responses to tune the model

Reinforcement learning from human feedback (RLHF) [PREVIEW](#)

Uses human preference data to create a separate reward model, which then tunes the foundation model using reinforcement learning

[CONTINUE](#)

Model name \*

The name of the new model. Up to 128 characters.

## Tuning settings

Base model

text-bison@001

The base model that will be used to create a new tuned model.

Train steps \*

300

Learning rate multiplier \*

1

The real learning rate is defined as the product of the learning rate multiplier and the underlying recommended learning rate.

 gs:// Working directory \*

BROWSE

The Cloud Storage location where the artifacts are stored during the pipeline tuning run.

Accelerator type

GPU (us-central1)

The accelerator used for model tuning. The accelerator type you choose determines the region where the model tuning computation occurs.

Region

us-central1 (Iowa)

Where your pipeline tuning job runs and the tuned model is deployed. Your tuning data always remains in this region

## Add a TensorBoard instance

TensorBoard instance

# Tune a Model

← Create a tuned model

Tuning type

Model details

**3**  Tuning dataset

**4**  Evaluation (optional)

START TUNING

## Tuning dataset

The [tuning dataset](#) is a JSONL file that contains model prompt and responses examples(one per line). It's recommended that you use at least [100-500 samples](#).

[View sample dataset](#)

You can upload the file or select one that's already on Cloud Storage.

Upload JSONL file to Cloud Storage

Existing JSONL file on Cloud Storage

Select JSONL file \*

BROWSE

Each line in the file contains one example:

- An `input_text` field containing the prompt
- An `output_text` field containing an example response

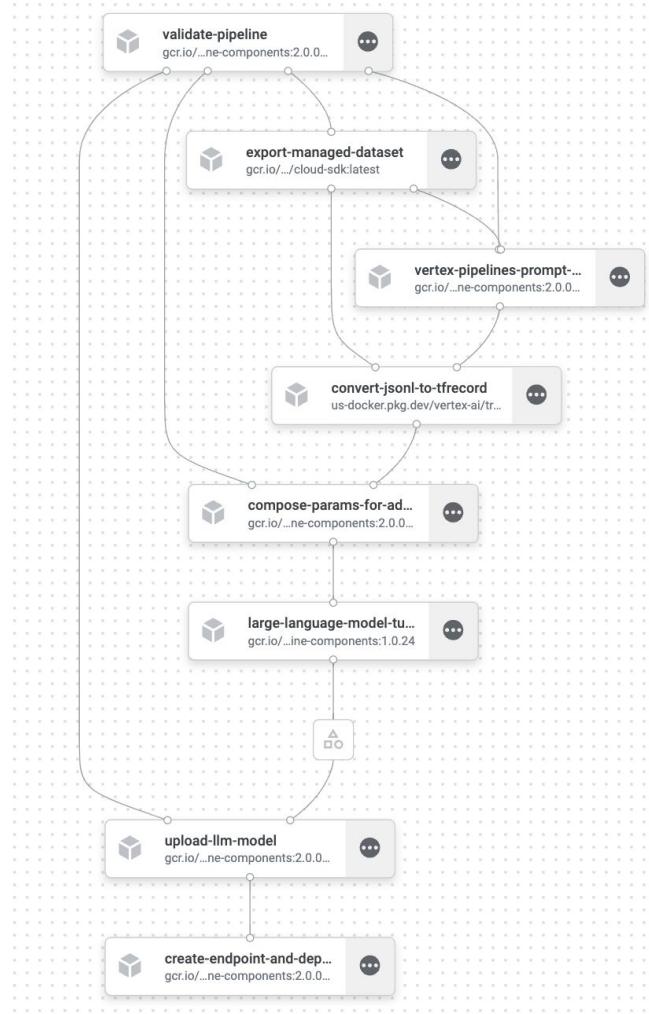
For example: `{"input_text": "Create a description for Plantation Palms", "output_text": "Enjoy some fun in the sun at Gulf Shores."}`

gs:// Dataset location \*

BROWSE

The Cloud Storage location where the JSONL file will be stored.

CONTINUE



# Model Garden

- 01 Introduction to Model Garden
- 02 Available Model Types
- 03 Working with Gen AI Studio
- 04 **Working with Model Registry**
- 05 Intro: Course Use Cases



# Model Registry

## Model Registry

[+ CREATE](#)[IMPORT](#)[REFRESH](#)[LEARN](#)

Models are built from your datasets or unmanaged data sources. There are many different types of machine learning models available on Vertex AI, depending on your use case and level of experience with machine learning. [Learn more](#)

Region

us-central1 (Iowa)

**Filter** Enter a property name

<input type="checkbox"/> Name	Deployment status	Description	Default version	Type	Source	Updated ↓	Labels
-------------------------------	-------------------	-------------	-----------------	------	--------	-----------	--------

No rows to display

Models are built from your datasets or unmanaged data sources. Choose the different types of machine learning models available on Vertex AI based on your use case and level of experience with machine learning. [Learn more](#)

## Region

us-central1 (Iowa)

**Filter** Enter a property name

<input type="checkbox"/> Name	Deployment status
No rows to display	

## Train new model

- 1 Training method
- 2 Model details
- 3 Training container
- 4 Hyperparameters (optional)
- 5 Compute and pricing
- 6 Prediction container (optional)

[START TRAINING](#)[CANCEL](#)

Dataset \*

No managed dataset



Annotation set

-



Objective

Custom



Please refer to the pricing guide for more details (and available deployment options) for each method.



AutoML options are only available when you train with a managed dataset.

### Model training method

 AutoML

Train high-quality models with minimal effort and machine learning expertise. Just specify how long you want to train. [Learn more](#)

 AutoML Edge

Train a model that can be exported for on-prem/on-device use. Typically has lower accuracy. [Learn more](#)

 Custom training (advanced)

Run your TensorFlow, scikit-learn, and XGBoost training applications in the cloud. Train with one of Google Cloud's pre-built containers or use your own. [Learn more](#)

[CONTINUE](#)

Models are built from your datasets or unmanaged data sources. There are different types of machine learning models available on Vertex AI, depending on the case and level of experience with machine learning. [Learn more](#)

## Region

us-central1 (Iowa)

**Filter** Enter a property name

<input type="checkbox"/>	Name	Deployment status
--------------------------	------	-------------------

No rows to display

## Import model

- 1 Name and region
- 2 Model settings
- 3 Explainability (optional)

**IMPORT**

CANCEL

You can import model artifacts that have been trained outside of Google Cloud. Once your model has been imported, you can serve it for online or batch predictions and compare it against your other Cloud AI models. [More info](#)

 Import as new model

Creates a new model group and assigns the imported model as version 1

 Import as new version

Imports the model as a version of an existing model

Name \*

Description

Region

us-central1 (Iowa)

**▼ ADVANCED OPTIONS****CONTINUE**

# Model Card

## Stable Diffusion v1-5

Latent text-to-image diffusion model capable of generating photo-realistic images given a text input.

[DEPLOY](#)[OPEN NOTEBOOK](#)[OVERVIEW](#)[USE CASES](#)[DOCUMENTATION](#)[PRICING](#)

### Overview

The Stable Diffusion model was first published in the paper "[High-Resolution Image Synthesis with Latent Diffusion Models](#)" by Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer in CVPR 2022. To enable diffusion models training on limited computational resources while retaining their quality and flexibility, the authors applied them in the latent space of powerful pretrained autoencoders. Training diffusion models on such a representation allows for the first time to reach a near-optimal point between complexity reduction and detail preservation, greatly boosting visual fidelity. The latent diffusion models achieved a new state of the art for image inpainting and highly competitive performance on various tasks, including unconditional image generation, semantic scene synthesis, and super-resolution, while significantly reducing computational requirements compared to pixel-based diffusion models.

Supported by Stability AI and LAION, the 1st version of Stable Diffusion was trained on 512x512 images from a subset of the LAION-5B database. The model uses a frozen CLIP ViT-L/14 text encoder to condition the model on text prompts. With its 860M UNet and 123M text encoder, the model is relatively lightweight and runs on a GPU with at least 10GB VRAM.

The Stable Diffusion V1.5 was trained by runwayml which published the modeling implementation in their Github repository [runwayml/stable-diffusion](#).

### Resource ID

runwayml/stable-diffusion-v1-5

### Tags

Task

Generation

# Stable Diffusion v1-5

Latent text-to-image diffusion model capable of generating images given a text input.

[OVERVIEW](#)[USE CASES](#)[DOCUMENTATION](#)

## Save model

To deploy this model, you will need to save a copy to Model Registry in your Cloud project.

Resource ID

runwayml/stable-diffusion-v1-5 (Latest model)



Save as new model

Creates a new model group and assigns this model as version 1

Save as new version

Saves the model as a version of an existing model

Model name \*

stable\_diffusion



Some locations have been restricted due to a policy set by your organization. [Learn more about restricting locations.](#)

Region \*

us-central1 (Iowa)



## Advanced Options



**SAVE**

CANCEL

## Overview

The **Stable Diffusion** model was first published in the paper ["Stable Diffusion: A Latent Space Model for Image Generation and Inpainting"](#) by Robin Rombach, Andreas Blattmann, Dominik Liebenwein, and Philipp Hennig. The authors applied them in the latent space of powerful pre-trained models. This representation allows for the first time to reach a near-optimal trade-off between visual fidelity preservation, greatly boosting visual fidelity. The latent diffusion model achieves state-of-the-art inpainting and highly competitive performance on various tasks such as image generation, scene synthesis, and super-resolution, while significantly reducing the computational cost compared to state-of-the-art diffusion models.

Supported by Stability AI and LAION, the 1st version of Stable Diffusion was trained on a dataset of 400M images from the LAION-5B database. The model uses a frozen CLIP VQGAN encoder to generate images from text prompts. With its 860M UNet and 123M text encoder, the model requires at least 10GB VRAM.

The Stable Diffusion V1.5 was trained by runwayml which provides a public repository [runwayml/stable-diffusion](#).

# Stable Diffusion v1-5

Latent text-to-image diffusion model capable of generating images given a text input.

**DEPLOY** OPEN NOTEBOOK [🔗](#)

[OVERVIEW](#) [USE CASES](#) [DOCUMENTATION](#)

## Overview

The **Stable Diffusion** model was first published in the paper "Stable Diffusion: A Latent Space Model for Image Generation and Inpainting" by Robin Rombach, Andreas Blattmann, Dominik Reichert, and Philipp Hennig. This model enables diffusion models training on limited computational resources. The authors applied them in the latent space of powerful pre-trained models. This representation allows for the first time to reach a near-optimal balance between preservation and efficiency. The latent diffusion model achieves state-of-the-art results in inpainting and highly competitive performance on various tasks such as scene synthesis, and super-resolution, while significantly outperforming other diffusion models.

Supported by Stability AI and LAION, the 1st version of Stable Diffusion uses a frozen CLIP model trained on 100M images from the LAION-5B database. The model uses a frozen CLIP model for text embeddings. With its 860M UNet and 123M text encoder, the model requires at least 10GB VRAM.

The Stable Diffusion V1.5 was trained by runwayml which has a public GitHub repository [runwayml/stable-diffusion](#).

**Deploy to endpoint**

Create new endpoint  Add to existing endpoint

Endpoint name \* [?](#)

**Location**

Some locations have been restricted due to a policy set by your organization. [Learn more about restricting locations.](#) [🔗](#)

Region: us-central1 (Iowa) [?](#)

**Access**

Determines how your endpoint can be accessed. By default, endpoints are available for prediction serving through a REST API. Endpoint access can't be changed after the endpoint is created.

Standard  
Makes the endpoint available for prediction serving through a REST API. AutoML and custom-trained models can be added to standard endpoints.

Private  
Create a private connection to this endpoint using a VPC network and [private services access](#). Only custom-trained and tabular models can be added to private endpoints. [Learn more](#) [🔗](#)

ADVANCED OPTIONS

**CONTINUE**

# Model versions

[stable\\_diffusion](#)[EDIT DETAILS](#)**Model description**

—

**Region**

us-central1 (Iowa)

**Model labels**

—

**Versions****Filter** Enter a property name

<input type="checkbox"/> Version ID ↓	Alias	Status	Description	Endpoints	Created	Labels
<input type="checkbox"/> 1	<span>default</span>	<span>Ready</span>	—	—	Jun 29, 2023, 7:48:55 AM	<span>⋮</span>

# Model Garden

- 01 Introduction to Model Garden
- 02 Available Model Types
- 03 Working with Gen AI Studio
- 04 Working with Model Registry
- 05 **Intro: Course Use Cases**



# Course Use Cases

03

## Task-Specific Models

Problem focus: text classification

Vertex AI provides multiple ways to solve classification tasks. Task-specific models are an easy way to get started.

04

## Foundation Models

Problem focus: language understanding

The foundation models provide a lot of power just as they are. We can use them to solve many tasks.

05

## Fine-tunable models

Problem focus:  
Productionalization

Vertex AI provides us the tools to tune models for our own use cases. We're going to see how to formalize that into a pipeline.

# Lab:

## Exploring Model Garden

