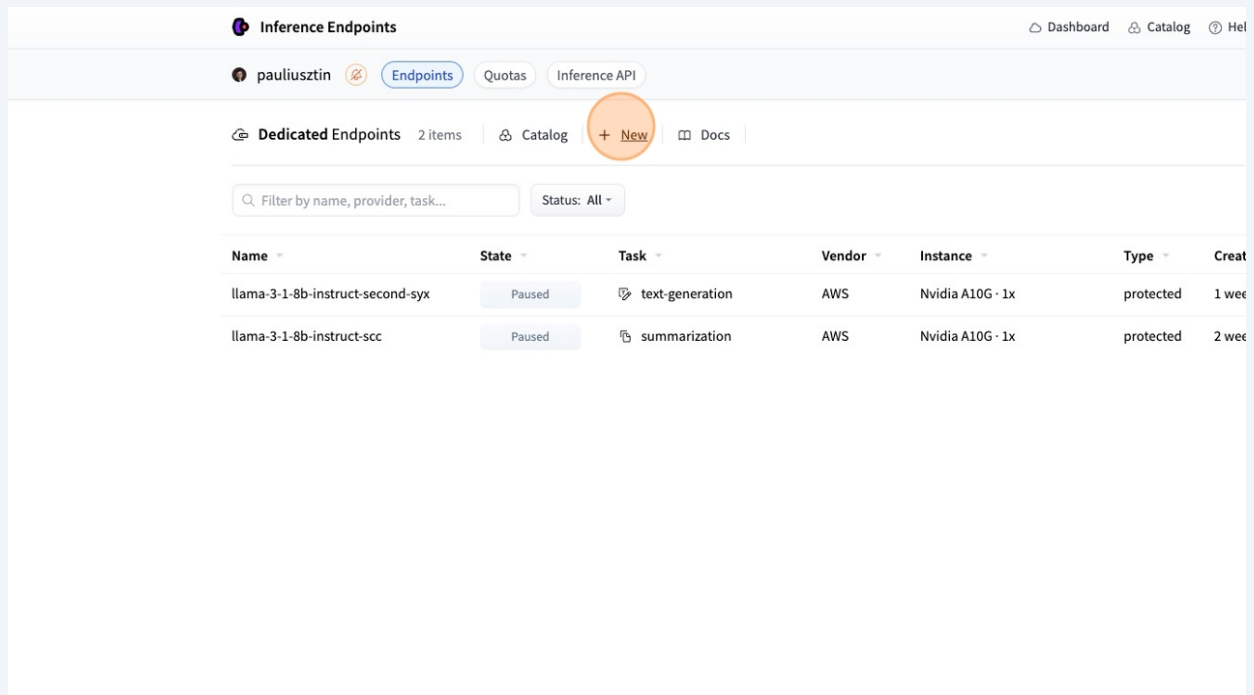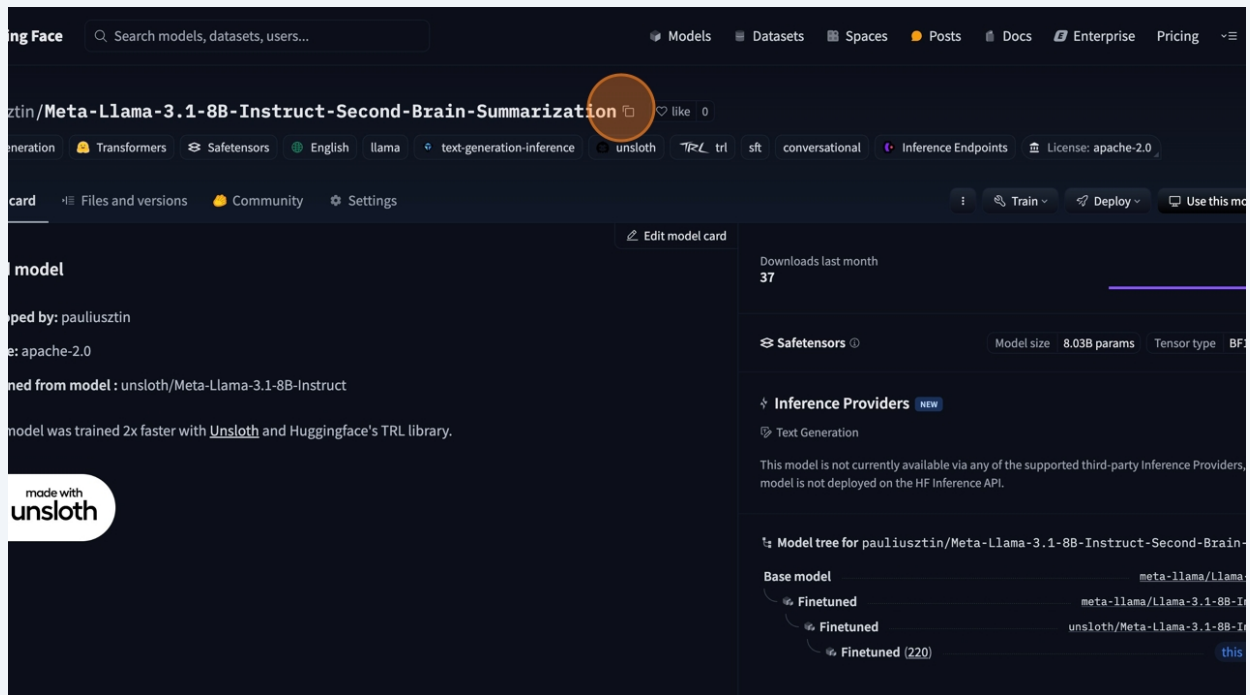# Creating an Inference Endpoint on Hugging Face

Scribe

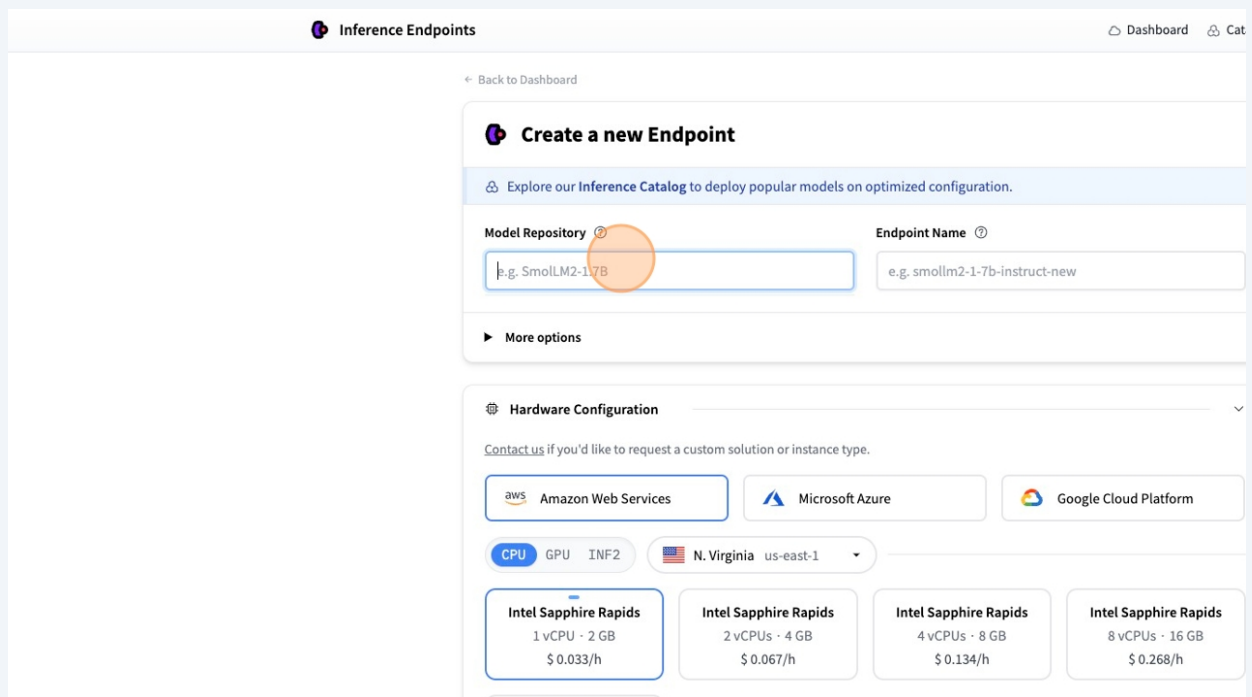**1**  Click "New"



**2**  Switch to tab
pauliusztin/Meta-Llama-3.1-8B-Instruct-Second-Brain-Summarization · Hugging Face"

**3**   Click this icon.



**4**   Switch to tab Create a new Endpoint | Inference Endpoints by Hugging Face"

**5**  Enter the Hugging Face model ID you copied in the "Model Repository" field.



**6**  We used "pauliusztin/Meta-Llama-3.1-8B-Instruct-Second-Brain-Summarization":

**7**    Click "GPU"

Model Repository ⓘ              Endpoint Name ⓘ

pauliusztin/Meta-Llama-3.1-8B-Instruct-Second-B... ✕     meta-llama-3-1-8b-instruct-s-xei

▶ **More options**

⊕ **Hardware Configuration**

Contact us if you'd like to request a custom solution or instance type.

| aws Amazon Web Services | ◭ Microsoft Azure | ☁ Google Cloud Platfor |
|---|---|---|

CPU  GPU  INF2    🇺🇸 N. Virginia  us-east-1  ▾

| Intel Sapphire Rapids | Intel Sapphire Rapids | Intel Sapphire Rapids | Intel Sapphire |
|---|---|---|---|
| 1 vCPU · 2 GB | 2 vCPUs · 4 GB | 4 vCPUs · 8 GB | 8 vCPUs · 16 |
| $ 0.033/h | $ 0.067/h | $ 0.134/h | $ 0.268/h |

**Intel Sapphire Rapids**
16 vCPUs · 32 GB
Reserved

ⓘ You may want to select a GPU accelerated instance to use the optimized Text Generation container.

---

**8**    Choose an Nvidia A10G GPU.

▶ **More options**

⊕ **Hardware Configuration**       ⌄

Contact us if you'd like to request a custom solution or instance type.

| aws Amazon Web Services | ◭ Microsoft Azure | ☁ Google Cloud Platform |
|---|---|---|

CPU  GPU  INF2    🇺🇸 N. Virginia  us-east-1  ▾

| Nvidia T4 | Nvidia L4 | Nvidia A10G | Nvidia L40S |
|---|---|---|---|
| 1 GPU · 16 GB | 1 GPU · 24 GB | 1 GPU · 24 GB | 1 GPU · 48 GB |
| 3 vCPUs · 15 GB | 7 vCPUs · 30 GB | 6 vCPUs · 30 GB | 7 vCPUs · 30 GB |
| $ 0.5/h | $ 0.8/h | $ 1/h | $ 1.8/h |
| Nvidia T4 | Nvidia L4 | Nvidia A100 | Nvidia A10G |
| 4 GPUs · 64 GB | 4 GPUs · 96 GB | 1 GPU · 80 GB | 4 GPUs · 96 GB |
| 46 vCPUs · 192 GB | 47 vCPUs · 185 GB | 11 vCPUs · 145 GB | 46 vCPUs · 186 GB |
| $ 3/h | $ 3.8/h | $ 4/h | $ 5/h |
| Nvidia A100 | Nvidia L40S | Nvidia A100 | Nvidia L40S |
| 2 GPUs · 160 GB | 4 GPUs · 192 GB | 4 GPUs · 320 GB | 8 GPUs · 384 GB |
| 22 vCPUs · 290 GB | 47 vCPUs · 380 GB | 44 vCPUs · 580 GB | 190 vCPUs · 1532 GB |
| $ 8/h | $ 8.3/h | $ 16/h | $ 23.5/h |
| Nvidia A100 | | | |
| 8 GPUs · 640 GB | | | |

**9** Select the closest region to you, such as "Ireland [eu-west-1]"

Model Repository ⓘ                          Endpoint Name ⓘ

pauliusztin/Meta-Llama-3.1-8B-Instruct-Second-B... ✕          meta-llama-3-1-8b-instruct-s-xei

▶ More options

⊕ Hardware Configuration                                                    ⌄

Contact us if you'd like to request a custom solution or instance type.

| aws  Amazon Web Services | ▲ Microsoft Azure | ☁ Google Cloud Platform |

CPU **GPU** INF2          🇺🇸 Ireland  eu-west-1   ▾

| Nvidia T4 | **Nvidia A10G** | Nvidia T4 | Nvidia A100 |
|---|---|---|---|
| 1 GPU · 16 GB | **1 GPU · 24 GB** | 4 GPUs · 64 GB | 1 GPU · 80 GB |
| 3 vCPUs · 15 GB | 6 vCPUs · 30 GB | 46 vCPUs · 192 GB | 11 vCPUs · 145 GB |
| $ 0.5/h | **$ 1/h** | $ 3/h | $ 4/h |

| Nvidia A10G | Nvidia A100 | Nvidia A100 | Nvidia A100 |
|---|---|---|---|
| 4 GPUs · 96 GB | 2 GPUs · 160 GB | 4 GPUs · 320 GB | 8 GPUs · 640 GB |
| 46 vCPUs · 186 GB | 22 vCPUs · 290 GB | 44 vCPUs · 580 GB | 88 vCPUs · 1160 GB |
| $ 5/h | $ 8/h | $ 16/h | $ 32/h |

⌄ **Autoscaling**   0 to 1 Replica / Scale-to-zero after 15 min                ⌃

---

**10** Go to the "Configuration" section

| $ 1/h | $ 3/h | $ 4/h |

| 0G | Nvidia A100 | Nvidia A100 | Nvidia A100 |
|---|---|---|---|
| 5 GB | 2 GPUs · 160 GB | 4 GPUs · 320 GB | 8 GPUs · 640 GB |
| 6 GB | 22 vCPUs · 290 GB | 44 vCPUs · 580 GB | 88 vCPUs · 1160 GB |
| | $ 8/h | $ 16/h | $ 32/h |

0 to 1 Replica / Scale-to-zero after 15 min                ⌃

                                                           ⌄

| ◯ Public | ◯ Private |
|---|---|

int is available from the Internet, secured with TLS/SSL and
gging Face Token for Authentication.

nfiguration   Text Generation Inference                     ⌃

Variables   No env variables defined                        ⌃

🗐 Create with cURL     **Create Endpoint**

**11** Select the "Bitsandbytes" quantization option

point is available from the Internet, secured with TLS/SSL and
igging Face Token for Authentication.

nfiguration ∨

r is the easiest way to deploy endpoints, and is very flexible thanks to custom Inference Handlers. You can also select a
for **Text-Generation** inference, or link your own **Custom** container.

Inference ∨

els ⑦                          Quantization ⑦

[                  ∨ ]          [ Bitsandbytes          ∨ ]

per Query) ⑦ optional          Max Number of Tokens (per Query) ⑦ optional

[ lt             ]             [ Container default              ]

okens ⑦ optional              Max Batch Total Tokens ⑦ optional

[ lt             ]             [ Container default              ]

Variables    No env variables defined ∧

---

**12** Click "Create Endpoint"

igging Face Token for Authentication.

nfiguration ∨

r is the easiest way to deploy endpoints, and is very flexible thanks to custom Inference Handlers. You can also select a
for **Text-Generation** inference, or link your own **Custom** container.

Inference ∨

els ⑦                          Quantization ⑦

[                  ∨ ]          [ Bitsandbytes          ∨ ]

per Query) ⑦ optional          Max Number of Tokens (per Query) ⑦ optional

[ lt             ]             [ Container default              ]

okens ⑦ optional              Max Batch Total Tokens ⑦ optional

[ lt             ]             [ Container default              ]

Variables    No env variables defined ∧

⎙ Create with cURL        [ Create Endpoint ]

**13** Click "Notify me!"