```
Welcome to Ubuntu 22.04.1 LTS (GNU/Linux 5.15.153.1-microsoft-standard-WSL2 x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:     https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage


This message is shown once a day. To disable it please create the
/root/.hushlogin file.
root@hp:~# docker model -h
Flag shorthand -h has been deprecated, use --help
Usage:  docker model COMMAND

Docker Model Runner (EXPERIMENTAL)

EXPERIMENTAL:
  docker model is an experimental feature.
  Experimental features provide early access to product functionality. These
  features may change between releases without warning, or can be removed from a
  future release. Learn more about experimental features in our documentation:
  https://docs.docker.com/go/experimental/

Commands:
  df               Show Docker Model Runner disk usage
  inspect          Display detailed information on one model
  install-runner   Install Docker Model Runner (Docker Engine only)
  list             List the models pulled to your local environment
  logs             Fetch the Docker Model Runner logs
  package          Package a GGUF file into a Docker model OCI artifact, with optional licenses, and pushes it to the specified registry
  ps               List running models
  pull             Pull a model from Docker Hub or HuggingFace to your local environment
  push             Push a model to Docker Hub
  rm               Remove local models downloaded from Docker Hub
  run              Run a model and interact with it using a submitted prompt or chat mode
  status           Check if the Docker Model Runner is running
  tag              Tag a model
  uninstall-runner Uninstall Docker Model Runner
  unload           Unload running models
  version          Show the Docker Model Runner version

Run 'docker model COMMAND --help' for more information on a command.
```

```
root@hp:~# docker model pull ai/smollm2
Using cached model: 258.06 MB
root@hp:~# docker model list
MODEL NAME    PARAMETERS  QUANTIZATION    ARCHITECTURE   MODEL ID     CREATED       SIZE
ai/smollm2    361.82 M    IQ2_XXS/Q4_K_M  llama          354bf30d0aa3 3 months ago  256.35 MiB
root@hp:~# docker model run ai/smollm2 "what is Kubernetes?"
Kubernetes (K8s) is a container orchestration platform that allows you to manage and scale containers across a large cluster of nodes. It's a powerful tool
that helps you to keep your applications running efficiently, whether they're on-prem or in the cloud.

Think of Kubernetes as a chef, but instead of cooking, it cooks containers, and instead of a kitchen, it cooks the containers. It manages the containers, ma
kes them available to the cluster, and schedules their execution.

Kubernetes can handle any number of containers, and it's not limited to just Kubernetes. It's a robust and flexible platform that can be used for a wide ran
ge of use cases, from simple deployments of software applications to complex enterprise-level automation.

You can use Kubernetes to manage production environments, virtual machines, and even cloud-based applications. It's a great tool for anyone looking to moder
nize their infrastructure or take control of their infrastructure from the ground up.

In essence, Kubernetes is a container orchestration platform that simplifies the process of managing containers across a cluster, making it easier to mainta
in and scale your applications.
root@hp:~# docker model run ai/smollm2
Interactive chat mode started. Type '/bye' to exit.
> What is AI?
AI stands for Artificial Intelligence. It's a type of computer system that can learn, think, and solve problems on its own, much like humans. AI systems use
 algorithms and data to recognize patterns and make decisions. They're not as smart as humans, but they can perform tasks that typically require human intel
ligence, such as understanding speech, recognizing images, and making decisions.
> What is Kubernetes?
Kubernetes (K8s) is a container orchestration software system used for managing containers on a cluster of servers. It is a container orchestration tool tha
t automates the deployment, scaling, and management of containers. Kubernetes was created by Google to optimize the deployment and management of containers
on large-scale environments. It uses a cluster of pods to manage containers and their resources.

Kubernetes helps organizations to manage multiple containers and their associated resources in a centralized manner, without having to manage and orchestrat
e each container individually. This helps in reducing the overhead associated with container management and allows the organization to scale its container e
nvironment more efficiently.

Kubernetes has become an essential tool for containerization in the industry, and it has been widely adopted by large-scale organizations, including Google,
 Amazon, and Microsoft.
> root@hp:~git clone https://github.com/docker/hello-genai.git
fatal: destination path 'hello-genai' already exists and is not an empty directory.
root@hp:~# cd hello-genai
root@hp:~/hello-genai# ls
```

```
Dockerfile  LICENSE  README.md  docker-compose.yml  go-genai  node-genai  py-genai  run.sh  rust-genai
root@hp:~/hello-genai# ls -ltra
total 64
-rw-r--r-- 1 root root  1848 Jul 19 13:02 docker-compose.yml
-rw-r--r-- 1 root root 11357 Jul 19 13:02 LICENSE
-rw-r--r-- 1 root root   413 Jul 19 13:02 Dockerfile
-rw-r--r-- 1 root root    64 Jul 19 13:02 .gitignore
drwxr-xr-x 3 root root  4096 Jul 19 13:02 node-genai
drwxr-xr-x 4 root root  4096 Jul 19 13:02 go-genai
drwxr-xr-x 5 root root  4096 Jul 19 13:02 rust-genai
-rwxr-xr-x 1 root root   517 Jul 19 13:02 run.sh
drwxr-xr-x 4 root root  4096 Jul 19 13:02 py-genai
-rw-r--r-- 1 root root   170 Jul 19 13:04 .env
drwxr-xr-x 7 root root  4096 Jul 19 13:04 .
-rw-r--r-- 1 root root  1201 Jul 19 13:43 README.md
drwxr-xr-x 8 root root  4096 Jul 19 13:43 .git
drwx------ 6 root root  4096 Jul 19 13:46 ..
root@hp:~/hello-genai# vim .env
root@hp:~/hello-genai# ./run.sh
Using LLM model: ai/smollm2
Pulling Docker model...
Using cached model: 258.06 MB
Running Docker Compose...
[+] Building 4.9s (78/78) FINISHED
 => [internal] load local bake definitions                                      0.0s
 => => reading from stdin 1.23kB                                                 0.0s
 => [rust-genai internal] load build definition from Dockerfile                 0.1s
 => => transferring dockerfile: 860B                                            0.0s
 => [python-genai internal] load build definition from Dockerfile               0.1s
 => => transferring dockerfile: 1.10kB                                          0.0s
 => [node-genai internal] load build definition from Dockerfile                 0.1s
 => => transferring dockerfile: 329B                                            0.0s
 => [go-genai internal] load build definition from Dockerfile                   0.1s
 => => transferring dockerfile: 999B                                            0.0s
 => [go-genai internal] load metadata for docker.io/library/alpine:3.18         3.7s
 => [go-genai internal] load metadata for docker.io/library/golang:1.20-alpine  3.7s
 => [rust-genai internal] load metadata for docker.io/library/debian:bookworm-slim  3.7s
 => [rust-genai internal] load metadata for docker.io/library/rust:1.82         3.7s
 => [python-genai internal] load metadata for docker.io/library/python:3.11-slim  3.7s
 => [node-genai internal] load metadata for docker.io/library/node:20-alpine    3.7s
 => [auth] library/node:pull token for registry-1.docker.io                     0.0s
```

```
=> [auth] library/node:pull token for registry-1.docker.io                                              0.0s
=> [auth] library/debian:pull token for registry-1.docker.io                                            0.0s
=> [auth] library/alpine:pull token for registry-1.docker.io                                            0.0s
=> [auth] library/rust:pull token for registry-1.docker.io                                              0.0s
=> [auth] library/golang:pull token for registry-1.docker.io                                            0.0s
=> [auth] library/python:pull token for registry-1.docker.io                                            0.0s
=> [node-genai internal] load .dockerignore                                                             0.0s
=> => transferring context: 2B                                                                          0.0s
=> [python-genai internal] load .dockerignore                                                           0.0s
=> => transferring context: 271B                                                                        0.0s
=> [rust-genai internal] load .dockerignore                                                             0.1s
=> => transferring context: 2B                                                                          0.0s
=> [go-genai internal] load .dockerignore                                                               0.1s
=> => transferring context: 239B                                                                        0.0s
=> [node-genai 1/6] FROM docker.io/library/node:20-alpine@sha256:df02558528d3d3d0d621f112e232611aecfee7cbc654f6b  0.0s
=> [node-genai internal] load build context                                                             0.0s
=> => transferring context: 151B                                                                        0.0s
=> [python-genai builder 1/4] FROM docker.io/library/python:3.11-slim@sha256:139020233cc412efe4c8135b0efe1c7569d  0.0s
=> [python-genai internal] load build context                                                           0.1s
=> => transferring context: 590B                                                                        0.0s
=> [rust-genai internal] load build context                                                             0.0s
=> => transferring context: 718B                                                                        0.0s
=> [rust-genai stage-1 1/7] FROM docker.io/library/debian:bookworm-slim@sha256:6ac2c08566499cc2415926653cf2ed7c3  0.0s
=> [rust-genai builder 1/7] FROM docker.io/library/rust:1.82@sha256:d9c3c6f1264a547d84560e06ffd79ed7a799ce0bff09  0.0s
=> [go-genai builder 1/7] FROM docker.io/library/golang:1.20-alpine@sha256:e47f121850f4e276b2b210c56df3fda919127  0.0s
=> [go-genai stage-1 1/7] FROM docker.io/library/alpine:3.18@sha256:de0eb0b3f2a47ba1eb89389859a9bd88b28e82f5826b  0.0s
=> [go-genai internal] load build context                                                               0.0s
=> => transferring context: 627B                                                                        0.0s
=> CACHED [node-genai 2/6] WORKDIR /app                                                                  0.0s
=> CACHED [node-genai 3/6] COPY package*.json ./                                                         0.0s
=> CACHED [node-genai 4/6] RUN npm install                                                               0.0s
=> CACHED [node-genai 5/6] COPY . .                                                                      0.0s
=> CACHED [node-genai 6/6] RUN mkdir -p views                                                            0.0s
=> CACHED [rust-genai stage-1 2/7] WORKDIR /app                                                          0.0s
=> CACHED [rust-genai stage-1 3/7] RUN apt-get update && apt-get install -y curl ca-certificates && rm -rf /var/  0.0s
=> CACHED [rust-genai stage-1 4/7] RUN useradd -m nomadicmehul                                           0.0s
=> CACHED [rust-genai builder 2/7] WORKDIR /usr/src/rust-genai                                           0.0s
=> CACHED [rust-genai builder 3/7] RUN cargo new --bin rust-genai                                        0.0s
=> CACHED [rust-genai builder 4/7] WORKDIR /usr/src/rust-genai/rust-genai                                0.0s
=> CACHED [rust-genai builder 5/7] COPY Cargo.toml ./                                                    0.0s
=> CACHED [rust-genai builder 6/7] COPY src ./src                                                        0.0s
```

```
=> CACHED [rust-genai builder 7/7] RUN cargo build --release                                        0.0s
=> CACHED [rust-genai stage-1 5/7] COPY --from=builder /usr/src/rust-genai/rust-genai/target/release/rust-genai   0.0s
=> CACHED [rust-genai stage-1 6/7] COPY static/ ./static/                                           0.0s
=> CACHED [rust-genai stage-1 7/7] COPY templates/ ./templates/                                     0.0s
=> CACHED [python-genai builder 2/4] WORKDIR /app                                                   0.0s
=> CACHED [python-genai stage-1 3/9] RUN adduser --disabled-password --gecos "" nomadicmehul        0.0s
=> CACHED [python-genai builder 3/4] COPY requirements.txt .                                        0.0s
=> CACHED [python-genai builder 4/4] RUN pip install --no-cache-dir -r requirements.txt             0.0s
=> CACHED [python-genai stage-1 4/9] COPY --from=builder /usr/local/lib/python3.11/site-packages /usr/local/lib/   0.0s
=> CACHED [python-genai stage-1 5/9] COPY --from=builder /usr/local/bin /usr/local/bin              0.0s
=> CACHED [python-genai stage-1 6/9] RUN mkdir -p templates static/css static/js static/images && chown -R nomad   0.0s
=> CACHED [python-genai stage-1 7/9] COPY app.py requirements.txt ./                                0.0s
=> CACHED [python-genai stage-1 8/9] COPY templates/ templates/                                     0.0s
=> CACHED [python-genai stage-1 9/9] COPY static/ static/                                           0.0s
=> [rust-genai] exporting to image                                                                 0.1s
=> => exporting layers                                                                             0.0s
=> => writing image sha256:c180622abf766cd2add639ad10b905e5bd233bd0793a291583d9b7c1a69504a6        0.0s
=> => naming to docker.io/library/hello-genai-rust-genai                                           0.0s
=> [node-genai] exporting to image                                                                 0.1s
=> => exporting layers                                                                             0.0s
=> => writing image sha256:6e94f285e5edd2ad3e28d29362fa35928d266d15c8022354c79e9dcaf91b43df        0.0s
=> => naming to docker.io/library/hello-genai-node-genai                                           0.0s
=> [python-genai] exporting to image                                                               0.1s
=> => exporting layers                                                                             0.0s
=> => writing image sha256:b7edf2f1184b708037a5e7d15b872263a96a9a80e5b00c3b4a497567f3f69217        0.0s
=> => naming to docker.io/library/hello-genai-python-genai                                         0.0s
=> CACHED [go-genai stage-1 2/7] WORKDIR /app                                                       0.0s
=> CACHED [go-genai stage-1 3/7] RUN adduser -D -g '' nomadicmehul                                  0.0s
=> CACHED [go-genai builder 2/7] WORKDIR /app                                                       0.0s
=> CACHED [go-genai builder 3/7] RUN apk add --no-cache git                                         0.0s
=> CACHED [go-genai builder 4/7] COPY go.mod ./                                                     0.0s
=> CACHED [go-genai builder 5/7] RUN go mod tidy && go mod download                                 0.0s
=> CACHED [go-genai builder 6/7] COPY . .                                                           0.0s
=> CACHED [go-genai builder 7/7] RUN CGO_ENABLED=0 GOOS=linux go build -a -installsuffix cgo -o main .   0.0s
=> CACHED [go-genai stage-1 4/7] COPY --from=builder /app/main .                                    0.0s
=> CACHED [go-genai stage-1 5/7] COPY static/ ./static/                                             0.0s
=> CACHED [go-genai stage-1 6/7] COPY templates/ ./templates/                                       0.0s
=> CACHED [go-genai stage-1 7/7] RUN chown -R nomadicmehul:nomadicmehul /app                        0.0s
=> [go-genai] exporting to image                                                                   0.1s
=> => exporting layers                                                                             0.0s
=> => writing image sha256:5078885c6cf714940dd18942af004836af13a6601953cb10c415775482397f99        0.0s
```

```
 => => naming to docker.io/library/hello-genai-go-genai                                                          0.0s
 => [rust-genai] resolving provenance for metadata file                                                          0.1s
 => [node-genai] resolving provenance for metadata file                                                          0.1s
 => [go-genai] resolving provenance for metadata file                                                            0.0s
 => [python-genai] resolving provenance for metadata file                                                        0.0s
[+] Running 4/4
 ✔go-genai      Built                                                                                            0.0s
 ✔python-genai  Built                                                                                            0.0s
 ✔node-genai    Built                                                                                            0.0s
 ✔rust-genai    Built                                                                                            0.0s
Attaching to go-genai-1, node-genai-1, python-genai-1, rust-genai-1
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 Configuration loaded: Port=8080, LLM Base URL=http://model-runner.docker.internal/engines/llama.cpp/v1,
Model=ai/smollm2
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 Current working directory: /app
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 Static directory exists
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 swagger.json exists
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 Server starting on http://localhost:8080
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 Using LLM endpoint: http://model-runner.docker.internal/engines/llama.cpp/v1/chat/completions
go-genai-1     | [hello-genai] 2025/07/21 14:10:54 Using model: ai/smollm2
node-genai-1   | Server starting on http://localhost:8082
node-genai-1   | Using LLM endpoint: http://model-runner.docker.internal/engines/llama.cpp/v1/chat/completions
node-genai-1   | Using model: ai/smollm2
python-genai-1 | [2025-07-21 14:10:55,825] INFO in app: Server starting on http://localhost:8081
python-genai-1 | [2025-07-21 14:10:55,826] INFO in app: Using LLM endpoint: http://model-runner.docker.internal/engines/llama.cpp/v1/chat/completions
python-genai-1 | [2025-07-21 14:10:55,827] INFO in app: Using model: ai/smollm2
python-genai-1 |  * Serving Flask app 'app'
python-genai-1 |  * Debug mode: off
python-genai-1 | WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
python-genai-1 |  * Running on all addresses (0.0.0.0)
python-genai-1 |  * Running on http://127.0.0.1:8081
python-genai-1 |  * Running on http://172.20.0.4:8081
python-genai-1 | Press CTRL+C to quit
python-genai-1 | 172.20.0.1 - - [21/Jul/2025 14:11:26] "GET / HTTP/1.1" 200 -
python-genai-1 | 172.20.0.1 - - [21/Jul/2025 14:11:26] "POST /api/chat HTTP/1.1" 200 -
python-genai-1 | 172.20.0.1 - - [21/Jul/2025 14:11:26] "GET /static/favicon.ico HTTP/1.1" 304 -
python-genai-1 | 172.20.0.1 - - [21/Jul/2025 14:11:55] "POST /api/chat HTTP/1.1" 200 -
python-genai-1 | 172.20.0.1 - - [21/Jul/2025 14:13:30] "POST /api/chat HTTP/1.1" 200 -

v View in Docker Desktop    o View Config    w Enable Watch
```