

# TextSynth Server

Version: 2023-10-21

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Quick Start</b>	<b>2</b>
2.1	Linux	2
2.1.1	First steps	2
2.1.2	GPU usage	2
2.2	Windows	3
2.2.1	First steps	3
2.2.2	GPU usage	4
<b>3</b>	<b>Utilities</b>	<b>5</b>
3.1	Text processing ( <b>ts_test</b> )	5
3.2	Text to image ( <b>ts_sd</b> )	5
3.3	Chat ( <b>ts_chat</b> )	5
3.4	Text compression ( <b>ts_zip</b> )	5
3.5	Model Weight Conversion	6
3.6	Model Weight Quantization	6
<b>4</b>	<b>TS Server Configuration</b>	<b>7</b>
4.1	Syntax	7
4.2	Configuration parameters	7
4.3	JSON endpoints	8

# 1 Introduction

TextSynth Server is a web server proposing a REST API to large language models. They can be used for example for text completion, question answering, classification, chat, translation, image generation, ...

It has the following characteristics:

- All is included in a single binary. Very few external dependencies (Python is not needed) so installation is easy on Linux and Windows.
- Supports many Transformer variants (GPT-J, GPT-NeoX, OPT, Fairseq GPT, M2M100, CodeGen, GPT2, T5, RWKV, Llama 2, Falcon, MPT, Mistral) and Stable Diffusion.
- Integrated REST JSON API for text completion, translation and image generation.
- Integrated HTML GUI for testing.
- Very high performance for small and large batches on CPU and GPU. Support of dynamic batching to handle a large number of simultaneous requests.
- Efficient custom 8, 4 and 3 bit quantization.
- Larger models work optimally on lower cost GPUs (e.g. RTX 3090, RTX A6000) thanks to efficient quantization.
- Support of speculative sampling for even faster inference.
- Support of grammar based sampling to constraint the model output according to a BNF grammar or a JSON schema.
- Uses the LibNC library for simple tensor manipulation using the C language.
- Simple command line tools (`ts_test`, `ts_sd`, `ts_chat`, `ts_zip`) are provided to test the various models.

The free version is available only for non commercial use. Commercial organizations must buy the commercial version. The commercial version adds the following features:

- High performance handling of simultaneous requests.
- Memory usage limitation.
- Multiple models can be served at the same time.
- SSL/TLS requests.

## 2 Quick Start

### 2.1 Linux

#### 2.1.1 First steps

The TextSynth Server works only on x86 CPUs supporting AVX2 (all Intel CPUs since 2013 support it). The installation was tested on Fedora and CentOS/RockyLinux 8 distributions. Other distributions should work provided the `libjpeg` and `libmicrohttpd` libraries are installed.

1. Install the `libjpeg` and `libmicrohttpd` libraries. If you use Fedora, RHEL, CentOS or RockyLinux, you can type as root:

```
dnf install libjpeg libmicrohttpd
```

`ts_test` can be used without these libraries. `ts_sd` needs `libjpeg`. `ts_server` needs `libjpeg` and `libmicrohttpd`.

2. Extract the archive and go into its directory:

```
tar xtf ts_server-##version##.tar.gz
```

```
cd ts_server-##version##
```

when `##version##` is the version of the program.

3. Download one small example model such as `gpt2_117M.bin` from the `ts_server` web page.
4. Use it to generate text with the `"ts_test"` utility:

```
./ts_test -m gpt2-117M.bin g "The Linux kernel is"
```

The `-T` option can be used to use more or less CPU cores (the default is the number of physical cores).

5. Start the server:

```
./ts_server ts_server.cfg
```

You can edit the `ts_server.cfg` JSON configuration file if you want to use another model.

6. Try one request:

```
curl http://localhost:8080/v1/engines/gpt2_117M/completions \
-H "Content-Type: application/json" \
-d '{"prompt": "The Linux kernel is", "max_tokens": 100}'
```

The full request syntax is documented at <https://textsynth.com/documentation.html>.

7. You can use the integrated GUI by exploring with your browser:

```
http://localhost:8080
```

Now you are ready to load a larger model and to use it from your application.

#### 2.1.2 GPU usage

You need an Nvidia Ampere, ADA or Hopper GPU (e.g. RTX 3090, RTX 4090, RTX A6000, A100 or H100) in order to use the server with cuda 11.x or 12.x installed. Enough memory must be available to load the model.

1. First ensure that it is working on CPU (See [First steps], page 2).
2. Ensure that you have a compatible cuda installation with cuda 11.x or 12.x. The software is preconfigured with cuda 11.x. If you want to use cuda 12.x, then change the link to the `libnc_cuda.so` library:

```
ln -sf libnc_cuda-12.so libnc_cuda.so
```

3. Then try to use the GPU with the `ts_test` utility:

```
./ts_test --cuda -m gpt2-117M.bin g "The Linux kernel is"
```

If you get an error such as:

```
Could not load: ./libnc_cuda.so
```

it means that cuda is not properly installed or that there is a mismatch between the installed cuda version and the one `ts_server` was compiled with. You can use:

```
ldd ./libnc_cuda.so
```

to check that all the required cuda libraries are present on your system.

4. Then edit the `ts_server.cfg` configuration to enable GPU support by uncommenting

```
cuda: true
```

and run the server:

```
./ts_server ts_server.cfg
```

5. Assuming you have curl, Try one request:

```
curl http://localhost:8080/v1/engines/gpt2_117M/completions \
-H "Content-Type: application/json" \
-d '{"prompt": "The Linux kernel is", "max_tokens": 100}'
```

6. You can use the integrated GUI by exploring with your browser:

```
http://localhost:8080
```

7. Depending on the amount of memory available on your GPU, you can set the `memory` parameter in `ts_server.cfg` to limit the amount of memory used by the server. It is usually necessary to use a few gigabytes less than maximum available amount of GPU memory.

## 2.2 Windows

### 2.2.1 First steps

The TextSynth Server works only on x86 CPUs supporting AVX2 (all Intel CPUs since 2013 support it).

1. Extract the ZIP archive, launch the shell and go into its directory:

```
cd ts_server-##version##
```

when `##version##` is the version of the program.

2. Download one small example model such as `gpt2_117M.bin` from the `ts_server` web page.
3. Use it to generate text with the `"ts_test"` utility:

```
ts_test -m gpt2-117M.bin g "The Linux kernel is"
```

The `-T` option can be used to use more or less CPU cores (the default is the number of physical cores).

4. Start the server:

```
ts_server ts_server.cfg
```

You can edit the `ts_server.cfg` JSON configuration file if you want to use another model.

5. Assuming you installed curl (you can download it from <https://curl.se/windows/>), try one request:

```
curl http://localhost:8080/v1/engines/gpt2_117M/completions \
-H "Content-Type: application/json" \
-d '{"prompt": "The Linux kernel is", "max_tokens": 100}'
```

The full request syntax is documented at <https://textsynth.com/documentation.html>.

6. You can use the integrated GUI by exploring with your browser:

`http://localhost:8080`

Now you are ready to load a larger model and to use it from your application.

### 2.2.2 GPU usage

You need an Nvidia Ampere, ADA or Hopper GPU (e.g. RTX 3090, RTX 4090, RTX A6000, A100 or H100) in order to use the server with cuda 11.x or 12.x installed. Enough memory must be available to load the model.

1. First ensure that it is working on CPU (see the previous section).
2. Ensure that you have a compatible cuda installation with cuda 11.x or 12.x. The software automatically detects the cuda version.
3. Then try to use the GPU with the `ts_test` utility:

```
./ts_test --cuda -m gpt2-117M.bin g "The Linux kernel is"
```

If you get an error such as:

```
Could not load: libnc_cuda-12.dll (error=126)
```

it means that cuda is not properly installed.

4. Then edit the `ts_server.cfg` configuration to enable GPU support by uncommenting

```
cuda: true
```

and run the server:

```
./ts_server ts_server.cfg
```

5. Assuming you have curl, Try one request:

```
curl http://localhost:8080/v1/engines/gpt2_117M/completions \
-H "Content-Type: application/json" \
-d '{"prompt": "The Linux kernel is", "max_tokens": 100}'
```

6. You can use the integrated GUI by exploring with your browser:

`http://localhost:8080`

7. Depending on the amount of memory available on your GPU, you can set the `memory` parameter in `ts_server.cfg` to limit the amount of memory used by the server. It is usually necessary to use a few gigabytes less that maximum available amount of GPU memory.

## 3 Utilities

### 3.1 Text processing (ts\_test)

- Text generation

```
./ts_test --cuda -m gpt2_117M.bin g "Hello, my name is"
```

- Translation

```
./ts_test --cuda -m m2m100_1_2B_q8.bin translate en fr "The dispute \
focuses on the width of seats provided on long-haul flights for \
economy passengers."
```

assuming you downloaded the m2m100\_1\_2B\_q8.bin model.

- Short text compression and decompression

```
./ts_test --cuda -m gpt2_117M.bin cs "Hello, how are you ?"
```

```
./ts_test --cuda ds "##msg##"
```

where ##msg## is the compressed message.

When using a CPU, remove the --cuda option.

### 3.2 Text to image (ts\_sd)

```
./ts_sd --cuda -m sd_v1.4.bin -o out.jpg "an astronaut riding a horse"
```

assuming you downloaded sd\_v1.4.bin.

When using a CPU, remove the --cuda option.

### 3.3 Chat (ts\_chat)

```
./ts_chat --cuda -m llama2_7B_chat_q4.bin
```

assuming you downloaded llama2\_7B\_chat\_q4.bin.

When using a CPU, remove the --cuda option.

During the chat, some commands are available. Use /h during the chat to have some help. Type Ctrl-C once to stop the output and twice to quit.

### 3.4 Text compression (ts\_zip)

To compress a text file (here alice29.txt), assuming you downloaded the rwkv\_169M.bin model, use:

```
./ts_zip --cuda -m rwkv_169M.bin c alice29.txt /tmp/out.bin
```

To decompress it:

```
./ts_zip --cuda -m rwkv_169M.bin d /tmp/out.bin /tmp/out.txt
```

A checksum is included in the compressed file and it is automatically checked. It is essential to use the same software version, language model and GPU model when compressing and decompressing a file.

Large compression gains occur only if the input file is in a language that the language model has already seen.

The compression ratio, speed and memory usage depend on the language model but also on the selected context length (-l option) and batch size (-b option). They are both chosen automatically but can be overridden:

- The memory usage is proportional to the context length and batch size.

- The speed increases when the batch size increases.
- The compression ratio increases with larger context lengths and smaller batch sizes.

More information is available at [https://bellard.org/ts\\_server/ts\\_zip.html](https://bellard.org/ts_server/ts_zip.html).

### 3.5 Model Weight Conversion

TextSynth Server uses a specific file format to store the weights of the models. Python scripts are provided in `scripts/` to convert model checkpoints to the TextSynth format.

Example to convert the Pythia pytorch weights to TextSynth:

```
python gptneox_hf_convert.py config.json pytorch_model.bin \
    pythia_deduped_160M.bin
```

### 3.6 Model Weight Quantization

With the `ncconvert` utility, it is possible to quantize the model weights to 8, 4 or 3 bits. Quantization reduces the GPU memory usage and increases the inference speed. 8 bit quantization yields a negligible loss. 4 bit quantization yields a very small loss. 3 bit quantization currently only works on a GPU.

Examples:

8 bit quantization:

```
./ncconvert -q bf8 pythia_deduped_160M.bin pythia_deduped_160M_q8.bin
```

4 bit quantization:

```
./ncconvert -q bf4 pythia_deduped_160M.bin pythia_deduped_160M_q4.bin
```



## 4 TS Server Configuration

The file `ts_server.cfg` provides an example of configuration.

### 4.1 Syntax

The syntax is similar to JSON with a few modifications:

- property names can be unquoted  

```
{ property: 1 }
```
- Multi-line and single line C style comments are accepted

### 4.2 Configuration parameters

**cuda** Optional boolean (default = false). If true, CUDA (Nvidia GPU support) is enabled.

**device\_index** Optional integer (default = 0). Select the GPU device when using several GPUs. Use the `nvidia-smi` utility to list the available devices.

**n\_threads** Optional integer. When using a CPU, select the number of threads. It is set by default to the number of physical cores.

**full\_memory** Optional boolean (default = true). When using a GPU, `ts_server` reserves by default all the GPU memory for better efficiency. This parameter disables this behavior so that the GPU memory is allocated on demand.

**max\_memory** Optional integer (default = 0). If non zero, limit the consumed GPU memory to this value by pausing the HTTP requests until there is enough memory.  
 Since there is some overhead when handling the requests, it is better to set a value a few GB lower than the amount of total GPU memory.

**kv\_cache\_max\_count** Optional integer (default = 0). See the `kv_cache_size` parameter.

**kv\_cache\_size** Optional integer (default = 0). The KV cache is stored in CPU memory and is only used by the `chat` endpoint to store the context of the conversation to accelerate the inference. It is disabled by default. `kv_cache_size` sets the maximum KV cache memory in bytes. `kv_cache_max_count` sets the maximum number of cached conversations.

**models** Array of objects. Each element defines a model that is served. The following parameters are defined:

**name** String. Name (ID) of the model in the HTTP requests.

**filename** String. Filename of the model. You can use the conversion scripts to create one from Pytorch checkpoints if necessary.

**draft\_model** Optional string. Filename of a smaller model used to accelerate inference (speculative sampling). The draft model must use the same tokenizer as the large model.

`sps_k_max`

Optional integer. When using speculative sampling, specify the maximum number of tokens that is predicted by the draft model. The optimal value needs to be determined by experimentation. It is usually 3 or 4.

Note: the free version only accepts one model definition.

`local_port`

Integer. TCP port on which the HTTP server listens to.

`bind_addr`

Optional string (default = "0.0.0.0"). Set the IP address on which the server listens to. Use "127.0.0.1" if you want to accept local connections only.

`tls`

Optional boolean (default = false). If true, HTTPS (TLS) connections are accepted instead of HTTP ones.

`tls_cert_file`

Optional string. If TLS is enabled, the certificate (PEM format) must be provided with this parameter.

`tls_cert_file`

Optional string. If TLS is enabled, the private key of the certificate (PEM format) must be provided with this parameter.

`log_start`

Optional boolean (default = false). Print "Started." on the console when the server has loaded all the models and is ready to accept connections.

`gui`

Optional boolean (default = false). If true, enable a Graphical User Interface in addition to the remote API. It is available at the root URL, e.g. <http://127.0.0.1:8080>. The server just serves the files present in the `gui/` directory. You can modify or add new files if needed.

`log_filename`

String. Set the filename where the logs are written. There is one line per connection. The fields are:

- date and time (ISO format)
- source IP address
- HTTP method
- URI
- posted JSON

`from_proxy`

Optional boolean (default = true). If true, use the **X-Forwarded-For** header if available to determine the source IP address in the logs. It is useful to have the real IP address of a client when a proxy is used.

## 4.3 JSON endpoints

The server provides the following endpoints.

`v1/engines/{model_id}/completions`

Text completion.

Complete documentation at <https://textsynth.com/documentation.html>.

`v1/engines/{model_id}/chat`

Chat based completion. Complete documentation at <https://textsynth.com/documentation.html>.

`v1/engines/{model_id}/translate`

Translation.

Complete documentation at <https://textsynth.com/documentation.html>.

`v1/engines/{model_id}/logprob`

Log probability computation.

Complete documentation at <https://textsynth.com/documentation.html>.

`v1/engines/{model_id}/tokenize`

Tokenization.

Complete documentation at <https://textsynth.com/documentation.html>.

`v1/engines/{model_id}/text_to_image`

Text to image.

Complete documentation at <https://textsynth.com/documentation.html>.

`v1/memory_stats`

Return a JSON object with the memory usage statistics. The following properties are available:

`cur_memory`

Integer. Current used memory in bytes (CPU or GPU memory).

`max_memory`

Integer. Maximum used memory in bytes since the last call (CPU or GPU memory).

`kv_cache_count`

Integer. Number of entries in the KV cache count.

`kv_cache_size`

Integer. CPU Memory in bytes used by the KV cache.

`v1/models`

Return the list of available models and their capabilities. It is used by the GUI.