# Artificial intellegence stack
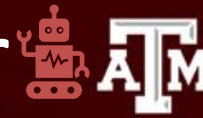
## *Stack*

# Required parts for a DL computer

**Motherboard:** Circuit board that connects all components together such as CPU, memory sticks etc., determines how compatible peripherals will attach themselves according to form factor standards such as ATX, Micro-ATX, Mini ITX etc...

**Central processing unit (CPU):** The brain of a computer; executes instructions from software programs and processes data; comes with different cores depending on how many threads may be active at once. Modern CPUs have multiple cores enabling multi-tasking capabilities. Must be able to handle tensor operations. **Intel Xeon, AMD Threadripper. 12C/24T

**Random Access Memory (RAM):** It refers to temporary memory used by a computer while it operates. Typically measured in megabytes or gigabytes (MB/GB). More RAM means faster performance and more programs can be open simultaneously. **128 Gb

**Storage:** Storage refers to retaining data on some form of device like SSD drives (Solid State Drives). Unlike traditional spinning disks known as Hard Disk Drives (HDDs), these newer alternative devices effectively provide faster read/write times by utilizing flash technology, often used also when retrieving files quickly. **4Tb SSD

**Graphics Processing Unit (GPU)/Tensor Processing Unit (TPU):** Processors designed specifically for handling graphics and computations related tasks like rendering images or performing deep learning operations respectively. NVIDIA is known for its powerful GPUs suitable for machine learning applications whereas Google developed TPUs which are optimized for their own AI workloads. *NVIDIA workstation cards (RTX A6000, A100 , or H100)
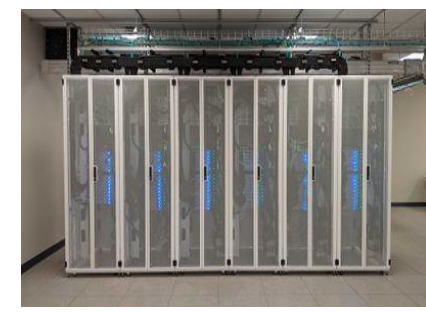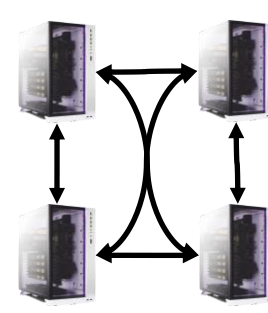
**Power Supply Unit (PSU):** Devices responsible for providing power to other hardware components within the system unit. They must match the requirements of the motherboard and other parts they power. Selecting an appropriate power supply unit guarantees smooth functioning if not exceeding wattage requirements could potentially cause instabilities or damage components within the system due to voltage fluctuations. **1100 watt

**Cooling System:** Important component preventing overheating caused by excessive heat generated during intensive activities – includes heatsinks/fans attached either manually or pre-installed upon purchase.- Considerations before purchasing include compatibility between each piece along with suitability concerning usage demands. For instance if one plans doing heavy artificial intelligence workstation activity then larger RAM and GPU will be used, which generate large amounts of heat. *6 fans @3000 RPM
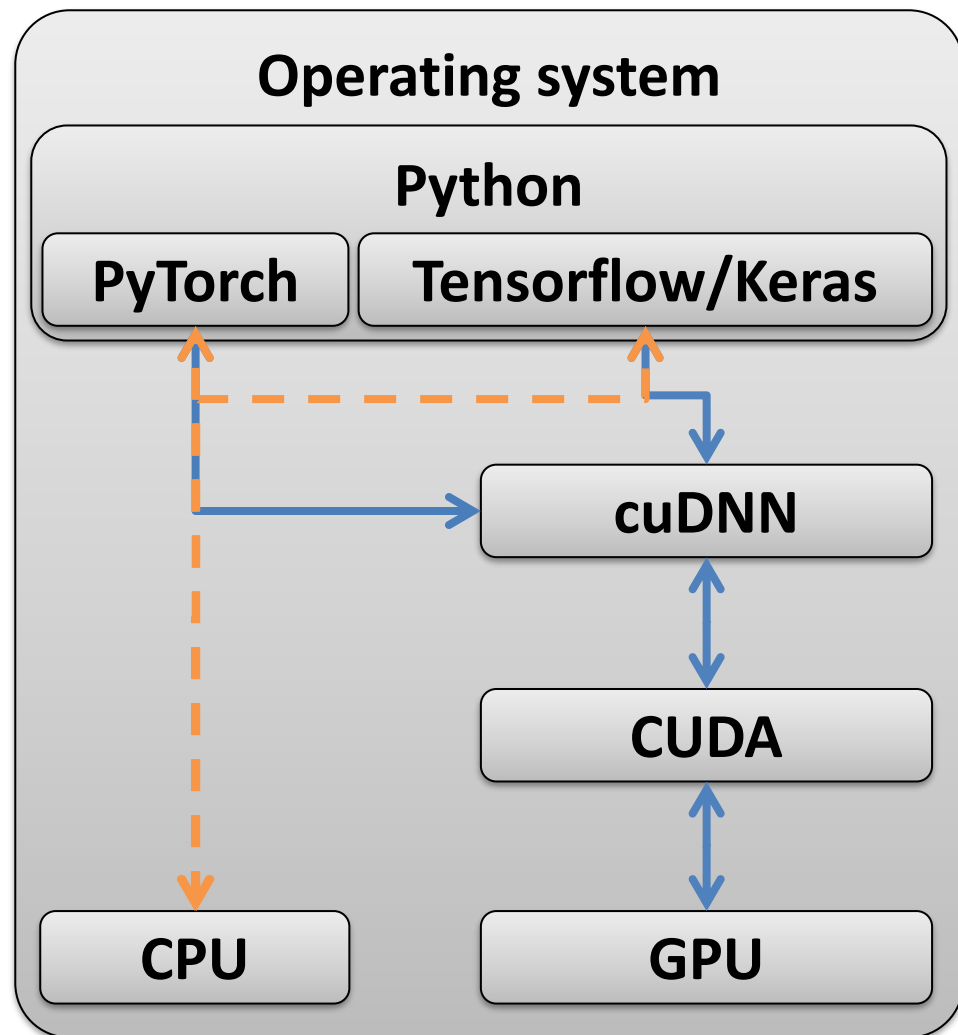
*Summarized with Llama 2 using Faraday

** Suggested minimum for a Ai/ML system

| | Laptop | Desktop | Workstations | Virtual Clusters | Centralized Clusters/HPC | Cloud |
|---|---|---|---|---|---|---|
| CPU | 4-10 Cores | 4-16 Cores | 12-64Cores | 100s | 1000s | " |
| GPU | 2-4 Gb | 1x-2x 16Gb | 2x-4x 48 Gb | 10s | 1000s | " |
| RAM | 8-32 Gb | 8 – 64 | 128 - 512 | 10s Tb | 100s Tb | " |
| Est. Cost | $1K-2K | $1K – 3K | $10K-35K | $50K> | $0.01 - 27.43/hr | " |
| Other fees/cost | Protection plan | Storage | Storage | Storage Software | Storage Software Infrastructure | Storage I/O Cont. Cost |
| Scaling options | Virtual cluster | Virtual cluster | Virtual cluster | Add macines | UTH, MDACC, UTMB, BCM, TAMU, Rice, Methodist, UH | $$$ |
| Access | Good | Good | Good | Administrated | Administrated | Open |

# Configuring a deep learning machine



| Version | Python version | Compiler | Build tools | cuDNN | CUDA |
|---|---|---|---|---|---|
| tensorflow-2.13.0 | 3.8-3.11 | Clang 16.0.0 | Bazel 5.3.0 | 8.6 | 11.8 |
| tensorflow-2.12.0 | 3.8-3.11 | GCC 9.3.1 | Bazel 5.3.0 | 8.6 | 11.8 |
| tensorflow-2.11.0 | 3.7-3.10 | GCC 9.3.1 | Bazel 5.3.0 | 8.1 | 11.2 |
| tensorflow-2.10.0 | 3.7-3.10 | GCC 9.3.1 | Bazel 5.1.1 | 8.1 | 11.2 |
| tensorflow-2.9.0 | 3.7-3.10 | GCC 9.3.1 | Bazel 5.0.0 | 8.1 | 11.2 |
| tensorflow-2.8.0 | 3.7-3.10 | GCC 7.3.1 | Bazel 4.2.1 | 8.1 | 11.2 |
| tensorflow-2.7.0 | 3.7-3.9 | GCC 7.3.1 | Bazel 3.7.2 | 8.1 | 11.2 |
| tensorflow-2.6.0 | 3.6-3.9 | GCC 7.3.1 | Bazel 3.7.2 | 8.1 | 11.2 |
| tensorflow-2.5.0 | 3.6-3.9 | GCC 7.3.1 | Bazel 3.7.2 | 8.1 | 11.2 |
| tensorflow-2.4.0 | 3.6-3.8 | GCC 7.3.1 | Bazel 3.1.0 | 8.0 | 11.0 |
| tensorflow-2.3.0 | 3.5-3.8 | GCC 7.3.1 | Bazel 3.1.0 | 7.6 | 10.1 |
| tensorflow-2.2.0 | 3.5-3.8 | GCC 7.3.1 | Bazel 2.0.0 | 7.6 | 10.1 |
| tensorflow-2.1.0 | 2.7, 3.5-3.7 | GCC 7.3.1 | Bazel 0.27.1 | 7.6 | 10.1 |
| tensorflow-2.0.0 | 2.7, 3.3-3.7 | GCC 7.3.1 | Bazel 0.26.1 | 7.4 | 10.0 |

tf: https://www.tensorflow.org/install/source#gpu

torch: https://pytorch.org/get-started/previous-versions/

# Choice of operating system

| OS: | Linux | Windows | WSL | MacOS |
|---|---|---|---|---|
| **Common:** | Ubuntu<br>CentOS<br>Red Hat | Windows 10<br>Windows 11 | WSL1<br>WSL2* | Ventura |
| **Command interface:** | Terminal (Bash) | Cmd, Powershell | Terminal (Bash) | Terminal (Bash) |
| **Pros:** | 1. Ubuntu is free<br>2. Open-source<br>3. Customizable<br>4. Less bloated<br>5. Common on HPRC | 1. Most widely used OS for professional systems<br>2. Large user base | 1. Lightweight Linux machine running on Windows 11<br>2. Best of both worlds<br>3. Supports Dockers | User interface<br>Linux-based |
| **Cons:** | 1. Step learning curve<br>2. Some compatibility issues<br>3. Commercial support issues | 1. Larger computational overhead<br>2. Less customizable/flexable | 1. Slightly reduced performance | Increased cost<br>Limited customizability |

**Pro's and con's generated by Wizard v1.1 model with GPT4ALL

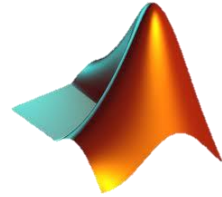| Host | Azure | AWS | Colab | OCI | DigitalOcean | IBM | Lambda |
|------|-------|-----|-------|-----|--------------|-----|--------|
| **Parent Co.** | Microsoft | Amazon | Google | Oracle | DigitalOcean, Inc. | IBM | Lambda Labs |
| **Interface** | Virtual compute Database Storage | Virtual compute Database Storage | Virtual Notebook | Virtual compute Database Storage | App hosting | Virtual compute Database Storage | Jupyter |

# Common code/scripting langs.

| Software: | Matlab | R statistics | Python |
|---|---|---|---|
| **Interface:** | Programing<br>Scripting | Programing<br>Scripting | Programing<br>Scripting |
| **Features:** | Machine Learning Toolbox | caret<br>e1071<br>randomForest<br>xgBoost<br>nueralnet | PIL<br>scikit-image<br>scikit-learn<br>Pytorch<br>Tensorflow |
| **Licensing:** | Commercial | Open Source | Open Source |

# Interactive Developer Environments (IDE)

| Software: | R studios | Jetbeans/PyCharm | Jupyter Labs/notebook |
|---|---|---|---|
| Languages: | R stats<br>Python | R stats<br>Python | R stats<br>Python |
| Features: | Scripting<br>Terminal<br>Notebooks<br>Env editor/preview | Scripting<br>Terminal<br>Env editor/preview | Scripting<br>Terminal<br>Notebooks<br>Env editor/preview |
| Licensing: | Free<br>Commercial | Free<br>Commercial | Free<br>Commercial |

# No-code & Low-code ML

| Software: | H2O Flow | Orange | Knime | Rapidminer | Pipeline Pilot |
|---|---|---|---|---|---|
| **Interface:** | Low-code | No-code | No-code<br>Low-code | No-code<br>Low-code | No-code<br>Low-code |
| **Features:** | AutoML<br>Notebook<br>GPT | Python with visual programing | Integration<br>Automation<br>Visualization | Automation<br>Visualization | Integration<br>Automation<br>Visualization |
| **Licensing:** | Open Source | Open Source | Open source | Commercial | Commercial |

# Chat-based AI interfaces



| Name | ChatGPT | GPT4ALL | Faraday.dev | H2Ogpt LLMstudios | Text-generation-webui | Jupyter_AI | Co-Pilot |
|------|---------|---------|-------------|-------------------|----------------------|------------|----------|
| **Parent Company** | OpenAI | GPT4ALL | Faraday | H2O.ai | oobabooga | Jupyter | Microsoft |
| **Release** | Alpha | Alpha | Alpha | Alpha | Alpha | Beta | Beta |
| **Base models** | GPT3.5 turbo GPT4 | HF | HF | HF | HF | GPT3.5 turbo | GPT3.5 turbo |
| **Instance** | Cloud | CPU_Local | CPU_Local GPU_Local | GPU_Local | GPU_Local | Cloud_API | Cloud |
| **Pricing model** | Free Subscription ($20/mo) | Free | Free | Free | Free | Free Requires API Keys | Subscription ($30/mo) |

# The new source of knowledge

## HuggingFace



## Github

# Huggingface models 🤗

| Name | GPT | Llama(2, 3) | Orca(2) | Wizard(Coder) | Falcon | T5 | BERT |
|---|---|---|---|---|---|---|---|
| Parent Co. | OpenAi | Meta | Microsoft | WizardLM | TII | Google | Google |
| Base Tech | GPT | Llama | Llama | Llama | Llama | t5 | BERT |
| Scope | Assistant Coding Text generation | Assistant Coding Text generation | Instruction Task Completion | Instruction Coding | Text generation | Seq2Seq | Text generation Text Classification NER |
| Sizes | | 7B 13B 70B | 7B 13B | 7B 13B 70B | 7B 40B 180B | | |