Author: **Sasa Kolovou (Athena RC)**. License: CC-BY (Creative Commons: Attribution)

Published at https://github.com/clarin-eric/standards on the 21st of March, 2025.

| ML Activity | Domain | File types |
|---|---|---|
| *Data collection, cleaning, tokenization* | **Data Preparation** | **Raw data formats:**<br>● **CSV**, **TSV**, **XLSX**: Tabular data.<br>● **JSON**, **XML**: Semi-structured data formats.<br>● **AVRO**, **PARQUET**, **ORC**: Optimized for big data frameworks like Hadoop and Spark (more specialized for big data systems and analytics platforms). |
| | | **Image/audio/video data:**<br>● **JPEG**, **PNG**, **BMP**: Images.<br>● **MP4**, **AVI**, **MKV**: Videos.<br>● **WAV**, **MP3**, **FLAC**: Audio files |
| | | **Text data:**<br>● **TXT**, **DOCX**, **PDF**: Plain and formatted text.<br>● **JSONL**: Line-delimited JSON, used for NLP datasets. |
| | | **Data annotations/labels:**<br>● **COCO JSON**, **PASCAL VOC XML**, **YOLO TXT**: Annotation formats for computer vision.<br>● **BIO**, **CONLL-U**: Annotation formats for NLP tasks. |
| *Pretraining/ fine-tuning, validation* | **Model Training** | **Configuration files:**<br>● **YAML**, **JSON**, **INI**: Hyperparameters, training configurations.<br>● **TOML**: A human-readable configuration format (.toml). |
| | | **Checkpoints and logs:**<br>● **HDF5**: Model weights (e.g., TensorFlow/Keras models).<br>● **PT**, **PTH**: PyTorch model checkpoints. |

| | | |
|---|---|---|
| | | • **PB**: Protocol Buffers for TensorFlow models.<br>• **ONNX**: Interchangeable format for trained models (.onnx file extension).<br>• **LOG**, **TXT**: Training logs and performance metrics.<br>• **TFRecord**: TensorFlow's format for serialized training examples (.tfrecord file extension).<br>• **Safetensors:** They are a file format for efficiently serialising and loading models with billions of parameters without the vulnerabilities found in pickle. |
| | | **Preprocessed datasets:**<br>• **NPY**, **NPZ**: NumPy array formats.<br>• **TFRecord**, **LMDB**: Serialized and database-like formats for efficiency.<br>• **Pickle (PKL)**: Serialized Python objects for data pipelines or model weights. |
| *Deployment, inference pipelines* | **Model Exchange** | **Model formats:**<br>• **ONNX**: Open Neural Network Exchange format for interoperability.<br>• **HDF5**, **PB**, **PT**, **PTH**, **JOBLIB**: Framework-specific formats for sharing models.<br>• **PMML**: Predictive Model Markup Language for statistical and machine learning models.<br>• **CoreML**: For Apple's ecosystem.<br>• **TFLite**: TensorFlow Lite for mobile and embedded devices.<br>• **TensorRT**: NVIDIA's format for optimized deployment. |
| | | **Containerized models:**<br>• **Docker images**: To encapsulate model-serving environments.<br>• **ZIP**, **TAR.GZ**: Bundles of models, configurations, and dependencies. |
| | | **Deployment configurations:**<br>• **YAML**, **JSON**: For APIs or container apps (e.g., Kubernetes, Docker Compose).<br>• **BentoML Bundles**, **MLflow models**: Framework-specific model packaging formats.<br>• **Serialized pipeline formats**: SKLearn pipelines saved as Pickle or Joblib files. |

------------------------------------------------------------------\*\*\*\*\*\*------------------------------------------------------------------

| LLM Activity | Domain | File types |
|---|---|---|
| *Data collection, cleaning, tokenization* | **Data Preparation** | **Raw data formats:**<br>• **CSV**, **TSV**, **XLSX**: Tabular data.<br>• **JSON**, **XML**: Semi-structured data formats.<br><br>**Preprocessed Data**:<br>• **Tokenized Data**: NPY, NPZ (NumPy arrays of tokenized sequences).<br>• **TFRecord**: Efficient format for large datasets.<br>• **Pickle (PKL)**: Serialized tokenized datasets (with caution for portability).<br><br>**Custom Formats**:<br>• Hugging Face datasets library supports Arrow (**.arrow**) and Parquet files (**.parquet**).<br><br>**Annotation Formats**:<br>• **JSONL** for NLP fine-tuning (e.g., prompt-response pairs for GPT-style models).<br>• **CONLL-U** for sequence labeling (NER, POS tagging). |
| *Pretraining/ fine-tuning, validation* | **Model Training** | **Configuration files:**<br>• **YAML**, **JSON**: Hyperparameters, tokenizer settings, training configurations.<br><br>**Intermediate Outputs**:<br>• **Model Weights**: PT, PTH, PB (PyTorch), HDF5 (TensorFlow/Keras/PyTorch).<br>• **Tokenizer States**: JSON files for BERT, GPT, and SentencePiece tokenizers.<br>• **Logs/Checkpoints**: TXT, CSV, TensorBoard logs, or MLflow files. |

| | | **Custom Datasets**: |
| --- | --- | --- |
| | | ● Hugging Face datasets (Arrow, JSONL, Parquet) or TFRecord. |
| *Deployment, inference pipelines* | **Model Exchange** | **Model formats:**<br>● **ONNX**: Open Neural Network Exchange format for interoperability.<br>● **PT**, **PTH, PB**: PyTorch-specific weights.<br>● **HF:** stands for Hugging Face's Transformers format. It is a binary format that stores the model's parameters in a compressed format.<br>● **HDF5**: .h5 or .hdf5 file extension<br>● **SavedModel**: TensorFlow weights (SavedModel is stored in a directory saved_model.pb finary file for the model, saved_model.pbtxt a human readable version of the same file,  and also we have two folders named *variables* and *assets*).<br>● **TFLite**: Optimized models for mobile (.tflite).<br>● **GGML (Glorot/Gated Gremlin MLmodel)**: Efficient, quantized model format for lightweight inference (e.g., for llama.cpp, gpt4all). Uses the **.bin** file extension.<br>● **GGUF (Glorot/Gated Gremlin Updatable Format):** updatable version of GGML that allows for fine-tuning or updating the model parameters. |
| | | **Tokenizer Files**:<br>● JSON, SentencePiece models (.model + .vocab). |
| | | **Metadata and Deployment Configs**:<br>● JSON, YAML for describing the model and tokenizer settings. |
| | | **Containerized models or Bundled Formats:**<br>● **Docker images**: To encapsulate model-serving environments.<br>● **MLflow Model Bundles** for deployment.<br>● Hugging Face Transformers repository format: Includes model (safetensors by default, model.bin), tokenizer (tokenizer.json), and configuration (config.json).<br>● BentoML format for LLM deployment. |