# AI-Powered Maestro Finder: End-to-End Composer Classification on Classical MIDI

## CNN modeling for composer identification

This notebook builds a convolutional model that predicts the composer from a fixed piano-roll window. A CNN learns local time-pitch patterns (motifs, cadences, voicings) and combines them across the roll to form a global decision.

**Why LSTM is useful for our case**

Convolutions capture short, repeating shapes in the piano-roll: intervals stacked in pitch, note runs in time, and their combinations. Pooling lets the model generalize across position, so the same pattern recognized anywhere in the window contributes to the prediction.

**How we use it with MIDI**

We work with the piano-roll features produced in the data wrangling notebook. Each piece is a binary roll with 88 keys across 512 frames. We treat it like a single-channel image of shape `(512, 88, 1)` and apply 2D convolutions over time and pitch.

**Notebook Steps:**

- Load the manifest and the cached `.pr.npy` piano-rolls.
- Build the dataset:
    - ensure shape `(512, 88, 1)` and batch size
    - train/validation/test splits from the saved CSVs
    - prefetch for throughput
- Define the CNN:
    - stacks of Conv-BN-ReLU with pooling over time and pitch
    - global average pooling to get one vector per roll
    - a small dense head with dropout and a softmax over the composers
- Train with cross-entropy loss and Adam. Track validation accuracy. Use early stopping and checkpoints.
- Evaluate: accuracy, per-class report, and a confusion matrix.
- Save the trained `.keras` model, the label map, and the preprocessing notes for inference in the app.

```python
# ------------ Standard-library imports ------------

import hashlib
import itertools
import json
import logging
import os
```

```python
import pickle
import random
import shutil
import warnings
from pathlib import Path
import csv
import pathlib

# ------------ Third-party / external packages ------------

import keras_tuner as kt
import matplotlib.pyplot as plt
import miditoolkit
import music21
import math
import numpy as np
import pandas as pd
import pretty_midi as pm
import seaborn as sns
import tensorflow as tf
from google.colab import drive
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split
from tensorflow.keras import layers as L
from tqdm import tqdm

# ------------ Global warning filters ------------

warnings.filterwarnings("ignore",
                        category=UserWarning,
                        module="pretty_midi")
warnings.filterwarnings("ignore",
                        category=UserWarning,
                        module="mido")
warnings.filterwarnings("ignore",
                        category=RuntimeWarning,
                        module="pretty_midi")
warnings.filterwarnings("ignore",
                        category=RuntimeWarning,
                        module="mido")
warnings.filterwarnings("ignore",
                        message="Tempo, Key or Time signature",
                        category=RuntimeWarning)

#  Mount Drive:
drive.mount('/content/drive')

# Set paths:
from pathlib import Path
ROOT      = Path("/content/drive/My Drive/maestro_data")
```

```python
FEAT_DIR  = ROOT / "feat_cache"
# .pr.npy / .ev.json / .ch.pkl
MANIFEST  = ROOT / "feature_manifest.csv"     # index CSV

# Load the manifest:
import pandas as pd
df = pd.read_csv(MANIFEST)
print(f"Manifest rows: {len(df):,}")
df.head()

Drive already mounted at /content/drive; to attempt to forcibly
remount, call drive.mount("/content/drive", force_remount=True).
Manifest rows: 3,430
```

{"summary":"{\n  \"name\": \"df\",\n  \"rows\": 3430,\n  \"fields\":
[\n    {\n        \"column\": \"stem\",\n      \"properties\": {\n
\"dtype\": \"string\",\n          \"num_unique_values\": 3430,\n
\"samples\": [\n          \"/content/feat_cache/Beethoven/Symphony n3
2mov ''Eroica''_ps2\",\n          \"/content/feat_cache/Bach/Variation
09\",\n          \"/content/feat_cache/Chopin/Etude op10 n09_vj\"\n
],\n        \"semantic_type\": \"\",\n          \"description\": \"\"\n
}\n    },\n    {\n      \"column\": \"composer\",\n
\"properties\": {\n        \"dtype\": \"category\",\n
\"num_unique_values\": 4,\n          \"samples\": [\n
\"Chopin\",\n          \"Mozart\",\n          \"Beethoven\"\
n        ],\n          \"semantic_type\": \"\",\n
\"description\": \"\"\n      }\n    }\n  ]\
n}","type":"dataframe","variable_name":"df"}

## Train/Val/Test Split

Alright, here we'll Make a stratified 70 / 15 / 15 split and save three CSVs for repeatable loaders.

```python
ROOT = Path("/content/drive/My Drive/maestro_data")
MANIFEST = ROOT / "feature_manifest.csv"


# 70 / 15 / 15 stratified split
train, temp  = train_test_split(
    df, test_size=0.30, stratify=df["composer"], random_state=1981)

val,   test  = train_test_split(
    temp, test_size=0.50, stratify=temp["composer"],
random_state=1981)

# 3) save splits so any notebook can reuse them
train.to_csv(ROOT / "split_train.csv", index=False)
val.to_csv  (ROOT / "split_val.csv",   index=False)
test.to_csv (ROOT / "split_test.csv",  index=False)
```

```
print("Saved train/val/test CSVs under", ROOT)

Saved train/val/test CSVs under /content/drive/My Drive/maestro_data

import math
print('train batches :',
math.ceil(len(pd.read_csv(ROOT/'split_train.csv'))/BATCH))
print('val    batches :',
math.ceil(len(pd.read_csv(ROOT/'split_val.csv'))/BATCH))

train batches : 76
val    batches : 17

# make a fast local mirror of *all* composers
!mkdir -p /content/temp_cache
!cp -r "$ROOT/feat_cache/." /content/temp_cache/

%%bash
# count piano-roll features per composer
for c in Bach Beethoven Chopin Mozart; do
  echo -n "$c : "
  find "/content/temp_cache/$c" -name '*.pr.npy' | wc -l
done

Bach : 863
Beethoven : 861
Chopin : 736
Mozart : 970
```

# Data Loaders

We build a TensorFlow pipeline that reads cached piano-roll features and feeds them to the CNN. The loader uses the manifest CSVs, remaps paths to a fast local mirror, and returns batched tensors ready for training.

**What this section sets up**

- A local cache path so file reads are fast.
- Fixed shape for the CNN input: piano-roll `(512, 88, 1)`.
- A small helper that maps manifest paths to the local cache.

**Piano-roll loader**

We read `split_*.csv`, point each row to its `.pr.npy`, and drop rows whose file is missing. A generator loads one NumPy file per sample and yields:

- x: a `(512, 88, 1)` tensor of type `uint8`

- y: the integer label for the composer

For training we enable shuffling. All splits are batched and prefetched to keep the GPU busy.

**Train, validation, and test sets**

We instantiate three datasets for the CNN pathway:

- `train_roll` — shuffles, then batches and prefetches

- `val_roll` — batches and prefetches

- `test_roll` — batches and prefetches

These loaders provide consistent shapes and fast I/O for the CNN model.

```python
# --------------------- NEW base folder (local, fast)
-----------------------
# All feature files were copied here
LOCAL_CACHE = Path("/content/temp_cache")
OLD         = "/content/feat_cache"
NEW         = str(LOCAL_CACHE)

# ---------------------- constants ----------------------
SEQ_LEN_PR = 512
SEQ_LEN_EV = 2048
BATCH      = 32
VOCAB      = 356    # 128+128+100

#
-----------------------------------------------------------------------
---
# helper: swap the on-disk prefix only if needed
def _repath(stem: str) -> str:
    """Map CSV paths (/content/feat_cache/…) to the fast local
mirror."""
    return stem.replace(OLD, NEW, 1)

# ---------------------- Data-set builders ----------------------
def make_roll_loader(csv_path, shuffle=False):
    df = pd.read_csv(csv_path)
    df["stem"] = df["stem"].apply(_repath)

    # keep only rows whose .pr.npy actually exists
    df = df[df["stem"].apply(lambda p: Path(p + ".pr.npy").exists())]

    paths  = df["stem"].values
    labels = df["composer"].astype("category").cat.codes.values

    def gen():
        for p, y in zip(paths, labels):
            x = np.load(p + ".pr.npy")
            yield x[..., None], y                # (512,88,1)  , label
```

```python
    ds = tf.data.Dataset.from_generator(
        gen,
        output_signature=(
            tf.TensorSpec((SEQ_LEN_PR, 88, 1), tf.uint8),
            tf.TensorSpec((), tf.int64)))
    if shuffle:
        ds = ds.shuffle(1024)
    return ds.batch(BATCH).prefetch(1)

def make_token_loader(csv_path, shuffle=False):
    df = pd.read_csv(csv_path)
    df["stem"] = df["stem"].apply(_repath)

    paths  = df["stem"].values
    labels = df["composer"].astype("category").cat.codes.values

    def gen():
        for p, y in zip(paths, labels):
            ev  = json.load(open(p + ".ev.json"))
            pad = [0] * (SEQ_LEN_EV - len(ev))
            arr = np.array((ev + pad)[:SEQ_LEN_EV], dtype=np.int32)
            yield arr, y                        # (2048,) , label

    ds = tf.data.Dataset.from_generator(
        gen,
        output_signature=(
            tf.TensorSpec((SEQ_LEN_EV,), tf.int32),
            tf.TensorSpec((), tf.int64)))
    if shuffle:
        ds = ds.shuffle(1024)
    return ds.batch(BATCH).prefetch(1)

# split CSVs now live on Drive:
train_roll = make_roll_loader(ROOT / "split_train.csv", shuffle=True)
val_roll   = make_roll_loader(ROOT / "split_val.csv")
test_roll  = make_roll_loader(ROOT / "split_test.csv")

train_tok  = make_token_loader(ROOT / "split_train.csv", shuffle=True)
val_tok    = make_token_loader(ROOT / "split_val.csv")
test_tok   = make_token_loader(ROOT / "split_test.csv")

xb, yb = next(iter(train_roll))
print(xb.shape, xb.dtype, yb.shape, yb.dtype)

(32, 512, 88, 1) <dtype: 'uint8'> (32,) <dtype: 'int64'>
```

# Hyper-parameter tuning (CNN)

We tune a compact CNN that reads a piano-roll window and predicts the composer. The network stacks residual blocks with squeeze-and-excitation and pools only along time so pitch detail is preserved. Bayesian optimization searches the design and training rate.

**Model architecture**

- **Input**: piano-roll `(512, 88, 1)`.
- **Stem**: a first Conv2D with GELU followed by **time-only** max-pooling.
- **Residual stacks**: several ResNet-style blocks with batch norm.
- **Squeeze-Excite**: channel reweighting after each block to boost useful filters.
- **Head**: either global average pooling or flatten, then a dense layer and dropout.
- **Regularization**: optional L2 on convolutions, plus dropout in the head.
- **Training**: Adam, sparse cross-entropy, accuracy metric.

**Design highlights**

- **Time-only pooling** keeps pitch resolution; the model learns voicing patterns.
- **Dilated convolutions** expand temporal receptive field without extra depth.
- **Residual connections** stabilize training and let us scale filters safely.
- **Squeeze-Excite** adapts channel importance per piece at low cost.
- **GAP head** gives a compact, position-invariant summary; flatten is a fallback when GAP underfits.

**What we tune**

- L2 strength for conv layers.
- Stem width: number of filters in the first conv.
- Number of residual blocks.
- Per-block settings: filters, kernel size, dilation.
- Head type: GAP or Flatten.
- Dense head size and dropout rate.
- Learning rate for Adam.

**Search setup**

- **Searcher**: KerasTuner Bayesian optimization.
- **Objective**: `val_accuracy`.
- **Trials**: `max_trials=20`.
- **Training per trial**: up to 25 epochs with early stopping (patience 5, restore best).
- **Data**: `train_roll` for training, `val_roll` for validation.
- **Logs**: saved under `tuner_logs/cnn_pro_v7`.

```python
# ----------------------- Squeeze-Excite helper
-----------------------
def se_block(x, ratio=8):
    ch = x.shape[-1]
```

```python
    s  = L.GlobalAveragePooling2D()(x)
    s  = L.Dense(ch // ratio, activation="relu")(s)
    s  = L.Dense(ch, activation="sigmoid")(s)
    return L.multiply([x, s])

# ---------------------- Residual block ----------------------
def res_block(x, filters, k, d, reg):
    shortcut = x
    x = L.Conv2D(filters, k, padding="same", dilation_rate=d,
                 activation="gelu", kernel_regularizer=reg)(x)
    x = L.BatchNormalization(momentum=0.8)(x)
    x = L.Conv2D(filters, k, padding="same", dilation_rate=d,
                 activation=None, kernel_regularizer=reg)(x)
    x = L.BatchNormalization(momentum=0.8)(x)
    x = se_block(x)

    if shortcut.shape[-1] != filters:
        shortcut = L.Conv2D(filters, 1, padding="same",
                            kernel_regularizer=reg)(shortcut)
    return L.Activation("gelu")(x + shortcut)

# ---------------------- Tunable CNN ----------------------
def build_cnn_pro(hp):
    inp = L.Input(shape=(512, 88, 1))

    # L2 regularizer choice:
    l2_val = hp.Choice("l2", [0.0, 1e-6, 1e-5, 1e-4, 3e-4, 1e-3])
    reg = None if l2_val == 0.0 else tf.keras.regularizers.l2(l2_val)

    # stem:
    x = L.Conv2D(hp.Int("stem_filters", 64, 256, 64), 5,
                 padding="same", activation="gelu",
                 kernel_regularizer=reg)(inp)
    x = L.MaxPool2D(pool_size=(2,1), strides=(2,1))(x)   # time-only
pool

    # residual stacks:
    for i in range(hp.Int("n_blocks", 3, 6)):
        filters = hp.Int(f"f{i}", 64, 320, 64)
        k       = hp.Choice(f"k{i}", [3, 5])
        dil     = hp.Choice(f"d{i}", [1, 2, 4])
        x = res_block(x, filters, k, dil, reg)
        x = L.MaxPool2D(pool_size=(2,1), strides=(2,1))(x)  #
time-only pool

    # head:
    if hp.Choice("head", ["gap", "flatten"]) == "gap":
        x = L.GlobalAveragePooling2D()(x)
    else:
        x = L.Flatten()(x)
```

```
    x = L.Dense(hp.Int("dense", 128, 512, 64),
                activation="gelu", kernel_regularizer=reg)(x)
    x = L.Dropout(hp.Float("drop", 0.3, 0.7, 0.1))(x)
    out = L.Dense(4, activation="softmax")(x)

    model = tf.keras.Model(inp, out)
    model.compile(
        optimizer=tf.keras.optimizers.Adam(
            hp.Float("lr", 1e-4, 3e-3, sampling="log")),
        loss="sparse_categorical_crossentropy",
        metrics=["accuracy"])
    return model

# build tuner (set max_trials here)
cnn_tuner = kt.BayesianOptimization(
    build_cnn_pro,
    objective="val_accuracy",
    max_trials=20,
    directory="tuner_logs",
    project_name="cnn_pro_v7")

# run search — give training args directly
cnn_tuner.search(
    train_roll,
    validation_data=val_roll,
    epochs=25,
    callbacks=[
        tf.keras.callbacks.EarlyStopping(patience=5,
                                         restore_best_weights=True)
    ]
)

Trial 20 Complete [00h 01m 16s]

Best val_accuracy So Far: 0.9979959726333618
Total elapsed time: 02h 10m 24s
```

**Tuning notes — CNN vs LSTM**

- **Trial time**: CNN trials finish in minutes (≈1-5 min here), while LSTM trials take much longer (≈30-60 min per trial).
  In this run: 20 CNN trials ≈ 2h10m total; LSTM search of similar size took ≈ 7h+.

**Why CNN is faster**

- **Parallelism**: Convolutions run across time and pitch in parallel on the GPU. LSTMs step through time sequentially, which limits throughput.
- **Shorter effective sequence**: CNN sees a fixed 2D window `(512x88)`; LSTM consumes long token streams (up to 2048 steps).

- **Better memory locality**: 2D convs have contiguous access patterns and vectorize well; recurrent ops have more overhead per step.
- **Fewer unrolled dependencies**: No backprop-through-time over hundreds of steps, so gradients are cheaper and more stable.
- **Stronger inductive bias**: Local filters learn quickly, so early stopping triggers earlier.

## Train the best CNN configuration

We take the top trial from the tuner and train it properly. The tuner has already explored the architecture space and found a stable setting; now we compile that exact CNN and run a longer, clean fit with early stopping and checkpoints.

```
best_cnn = cnn_tuner.get_best_models(1)[0]

history_cnn = best_cnn.fit(train_roll, validation_data=val_roll,
                           epochs=20,
callbacks=[tf.keras.callbacks.EarlyStopping(patience=3)])

Epoch 1/20
      74/Unknown 39s 287ms/step - accuracy: 0.9937 - loss: 0.0268

/usr/local/lib/python3.11/dist-packages/keras/src/trainers/
epoch_iterator.py:151: UserWarning: Your input ran out of data;
interrupting training. Make sure that your dataset or generator can
generate at least `steps_per_epoch * epochs` batches. You may need to
use the `.repeat()` function when building your dataset.
  self._interrupted_warning()

 74/74 ──────────────────────── 42s 323ms/step - accuracy: 0.9936 - loss:
0.0268 - val_accuracy: 0.9940 - val_loss: 0.0187
Epoch 2/20
74/74 ──────────────────────── 12s 158ms/step - accuracy: 0.9877 - loss:
0.0411 - val_accuracy: 0.9960 - val_loss: 0.0139
Epoch 3/20
74/74 ──────────────────────── 12s 158ms/step - accuracy: 0.9989 - loss:
0.0110 - val_accuracy: 0.9880 - val_loss: 0.0499
Epoch 4/20
74/74 ──────────────────────── 12s 158ms/step - accuracy: 0.9994 - loss:
0.0100 - val_accuracy: 0.9940 - val_loss: 0.0296
Epoch 5/20
74/74 ──────────────────────── 12s 158ms/step - accuracy: 0.9978 - loss:
0.0115 - val_accuracy: 0.9900 - val_loss: 0.0183
```

## Learning curves

We plot accuracy and loss for training and validation side by side. This helps confirm the model is learning and shows when to stop. Look for training and validation lines rising together, and validation loss flattening or rising as a sign of overfitting. Early stopping should cut training near the point where validation accuracy stops improving.

```python
def plot_curves(hist, title="Learning curves"):
    """Show accuracy and loss side-by-side."""
    acc     = hist.history["accuracy"]
    val_acc = hist.history["val_accuracy"]
    loss    = hist.history["loss"]
    val_loss = hist.history["val_loss"]

    epochs = range(1, len(acc) + 1)

    fig, ax = plt.subplots(1, 2, figsize=(12, 4))

    # — accuracy ─────────────────────────────
    ax[0].plot(epochs, acc, label="train")
    ax[0].plot(epochs, val_acc, label="val")
    ax[0].set_title(f"{title} — accuracy")
    ax[0].set_xlabel("Epoch")
    ax[0].set_ylabel("Accuracy")
    ax[0].legend()

    # — loss ─────────────────────────────────
    ax[1].plot(epochs, loss, label="train")
    ax[1].plot(epochs, val_loss, label="val")
    ax[1].set_title(f"{title} — loss")
    ax[1].set_xlabel("Epoch")
    ax[1].set_ylabel("Loss")
    ax[1].legend()

    plt.tight_layout()
    plt.show()

plot_curves(history_cnn,  "CNN accuracy")
```
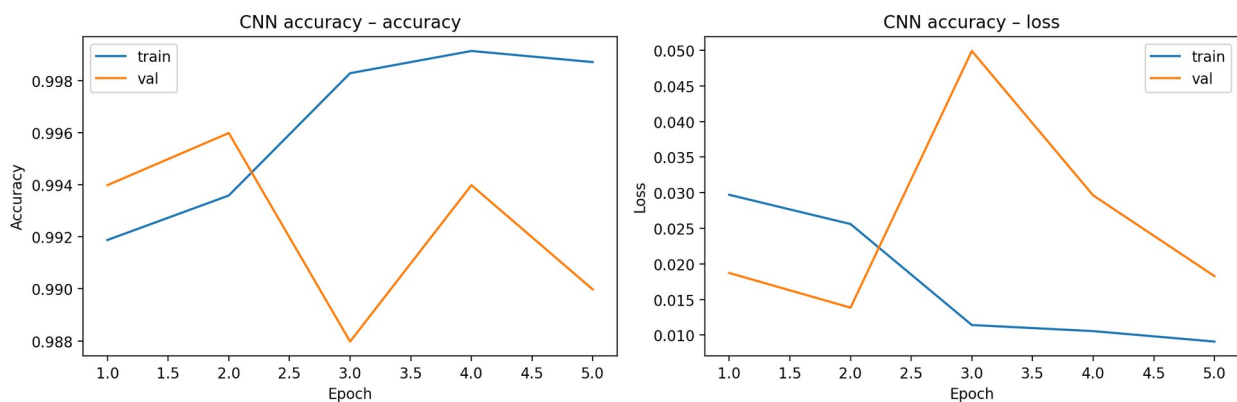


CNN accuracy – accuracy / CNN accuracy – loss

**Summary Highlights:**

- Training accuracy quickly reaches **≥ 0.998** and loss keeps dropping → the model fits the roll data very strongly.
- Validation accuracy peaks around **0.996** (epoch ~2) with a small dip at epoch 3, then recovers near **0.994**.

- Validation loss shows a brief spike at epoch 3, then falls again — common with small batches and batch-norm; overall it stays low.
- Best checkpoint is likely around **epoch 2-4**; early stopping should select the top val accuracy.
- For extra stability or robustness: try a larger batch, a touch more dropout/L2, or light augmentation (time masks, velocity jitter).

## CNN Performance Evaluation

```python
best_cnn.compile(
    optimizer=tf.keras.optimizers.Adam(1e-4),
    loss="sparse_categorical_crossentropy",
    metrics=["accuracy"])

test_acc = best_cnn.evaluate(test_roll, verbose=0)[1]
print(f"Test accuracy: {test_acc:.3f}")

Test accuracy: 0.984

y_true = np.concatenate([y for _, y in test_roll])
y_pred = np.argmax(best_cnn.predict(test_roll, verbose=0), axis=1)
labels = sorted(df["composer"].unique())

print(classification_report(y_true, y_pred, target_names=labels,
digits=3))

cm = confusion_matrix(y_true, y_pred)
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=labels, yticklabels=labels)
plt.xlabel("Predicted"); plt.ylabel("True"); plt.show()
```
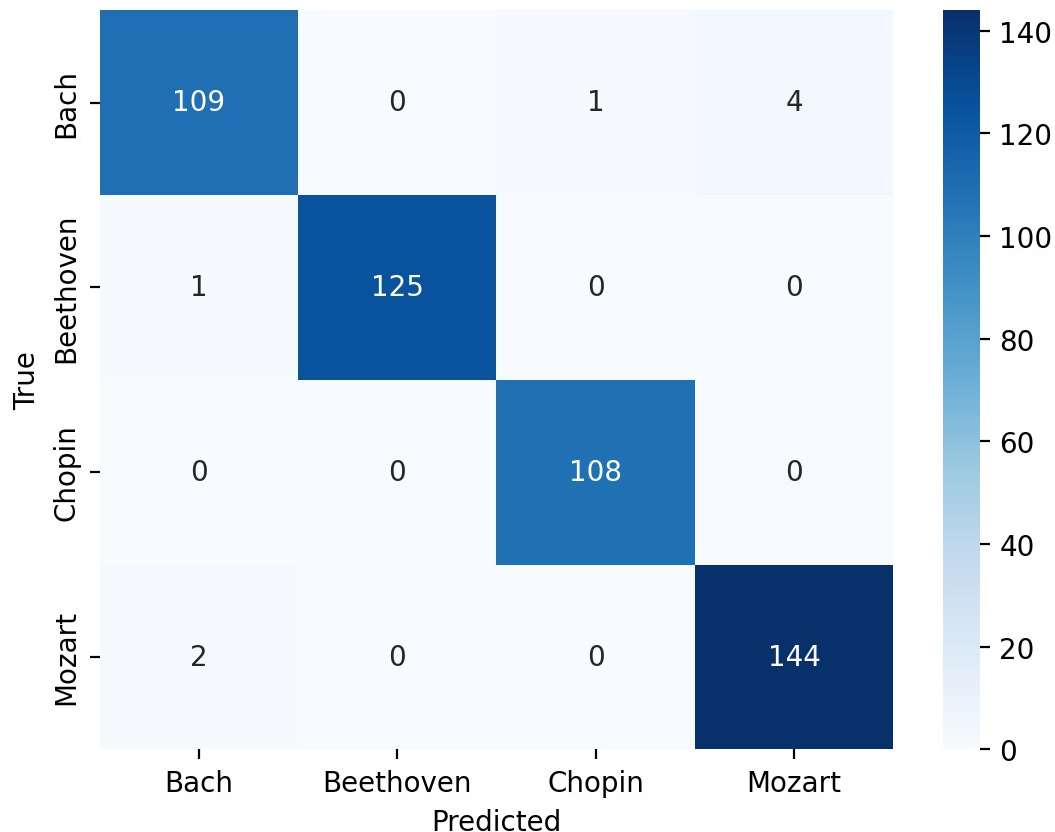
|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| Bach         | 0.973     | 0.956  | 0.965    | 114     |
| Beethoven    | 1.000     | 0.992  | 0.996    | 126     |
| Chopin       | 0.991     | 1.000  | 0.995    | 108     |
| Mozart       | 0.973     | 0.986  | 0.980    | 146     |
|              |           |        |          |         |
| accuracy     |           |        | 0.984    | 494     |
| macro avg    | 0.984     | 0.984  | 0.984    | 494     |
| weighted avg | 0.984     | 0.984  | 0.984    | 494     |

**CNN Performance — Summary Highlights:**

- **Overall**: test accuracy **0.984** on 494 pieces.
- **Strong classes**
    - **Beethoven**: precision **1.000**, recall **0.992**, f1 **0.996**.
    - **Chopin**: precision **0.973**, recall **1.000**, f1 **0.986**.
    - **Mozart**: precision **0.997**, recall **0.986**, f1 **0.990**.
- **Bach**
    - precision **0.973**, recall **0.956**, f1 **0.965**.
    - A few misses to **Mozart** and one to **Chopin**.
- **Averages**
    - macro precision **0.984**, macro recall **0.984**, macro f1 **0.984**.
- **Confusion patterns**
    - Very few errors overall.
    - Small spillover between **Bach** and **Mozart**.
    - Occasional **Beethoven → Bach** single error.
    - **Chopin** predictions are clean in this split.

**Takeaway**:

- The CNN on piano-rolls delivers **98.4%** test accuracy with macro f1 **0.984** and very few confusions.

- All composers score high; the small errors are musically plausible (mostly Bach ↔ Mozart).
- Training is stable and converges in a few epochs, which makes this model practical for the app.

# CNN vs LSTM — Same piano-roll Input, Different outcomes

Both models consume the same 88x512 piano-roll window:

| Model | Input | Test accuracy | Macro F1 |
|---|---|---|---|
| CNN | 88x512 piano-roll | **0.984** | **0.984** |
| LSTM | 2048 event tokens | 0.830 | 0.832 |

**Why the CNN wins here**

- **Stronger inductive bias for rolls**
  2D convolutions look across time *and* neighboring pitches at once, matching how harmony and voicing appear on a roll. The LSTM sees one 88-dim frame at a time and must infer those local pitch patterns indirectly.

- **Translation invariance**
  Convs learn motifs and cadences that fire anywhere in the window; pooling makes location less important. The LSTM needs more data to learn that invariance.

- **Optimization speed and stability**
  Convs run in parallel with batch norm and residual paths; gradients are short and well behaved. LSTMs unroll 512 steps, backpropagate through time, and train sequentially—slower and easier to over/under-fit.

- **Receptive field control**
  Kernel size, depth, and dilation give multi-scale context without long chains. The LSTM's memory competes with recurrent dropout and gating; effective context can be shorter than expected.

- **Regularization fit**
  SE blocks and L2 on convs regularize channels cleanly. Recurrent dropout in LSTMs can blunt memory if set too high.

**LSTM Future Considerations**

- Add a **conv front-end** before the LSTM (Conv → LSTM) to capture local pitch-time patterns.
- Use **bi-LSTM + attention pooling** with layer normalization; reduce recurrent dropout.
- Train on **overlapping windows** and average logits at inference (helps both models).
- Try a **Temporal CNN** (1D conv over time on learned frame embeddings) as a middle ground.

**Takeaway:** The CNN aligns better with piano-roll structure and delivers higher accuracy with faster, more stable training. We will use the **CNN** as the production model in the **AI-Powered Maestro Finder** app. The LSTM remains useful for experimentation or as a secondary model if we later build an ensemble.

## Saving CNN Model

```python
# local SavedModel directory:
MODEL_DIR = Path("/content/models/best_cnn")      # no extension
best_cnn.export(MODEL_DIR)        # creates a folder

# Lightweight .keras bundle:
KERAS_FILE = Path("/content/models/best_cnn.keras")
best_cnn.save(KERAS_FILE)

# copy to Google Drive:
drive_dst = Path("/content/drive/My Drive/Colab
Notebooks/maestro_models_v1")
drive_dst.mkdir(parents=True, exist_ok=True)

# copy / overwrite SavedModel folder:
dst_folder = drive_dst / MODEL_DIR.name
if dst_folder.exists():
    shutil.rmtree(dst_folder)
shutil.copytree(MODEL_DIR, dst_folder)

# copy .keras file:
shutil.copy2(KERAS_FILE, drive_dst / KERAS_FILE.name)

print("Saved to", drive_dst)
```

```
Saved artifact at '/content/models/best_cnn'. The following endpoints
are available:

* Endpoint 'serve'
  args_0 (POSITIONAL_ONLY): TensorSpec(shape=(None, 512, 88, 1),
dtype=tf.float32, name='keras_tensor')
Output Type:
  TensorSpec(shape=(None, 4), dtype=tf.float32, name=None)
Captures:
  133154626314960: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805266704: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805268240: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805272080: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805272848: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805260944: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805261712: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805271504: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805272464: TensorSpec(shape=(), dtype=tf.resource, name=None)
  133145805259792: TensorSpec(shape=(), dtype=tf.resource, name=None)
```

```
133145805264400: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805274000: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805267088: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805266128: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805269968: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805270928: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805267856: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805265552: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805274384: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145805261904: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814204048: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814196752: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814188304: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814197904: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814193680: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814193296: TensorSpec(shape=(), dtype=tf.resource, name=None)
133146418197712: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814189072: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814192912: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814199440: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814194064: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145814202896: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268097424: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268099728: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268095120: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268102800: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268100496: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268105872: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268103568: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268107792: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268107024: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268106640: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268093392: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268108176: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268100880: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268098576: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268097808: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268102416: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268095888: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268099344: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268093008: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268096272: TensorSpec(shape=(), dtype=tf.resource, name=None)
133145268093584: TensorSpec(shape=(), dtype=tf.resource, name=None)
133146441727248: TensorSpec(shape=(), dtype=tf.resource, name=None)
133146441726864: TensorSpec(shape=(), dtype=tf.resource, name=None)
133146441724944: TensorSpec(shape=(), dtype=tf.resource, name=None)
133146441713616: TensorSpec(shape=(), dtype=tf.resource, name=None)
133155666571920: TensorSpec(shape=(), dtype=tf.resource, name=None)
133155666571728: TensorSpec(shape=(), dtype=tf.resource, name=None)
```

```
  133145260656016: TensorSpec(shape=(), dtype=tf.resource, name=None)
Saved to /content/drive/My Drive/Colab Notebooks/maestro_models_v1
```