

Modelos MLP e CNN –
UrbanSound8K

Classificação de Áudio com Deep Learning

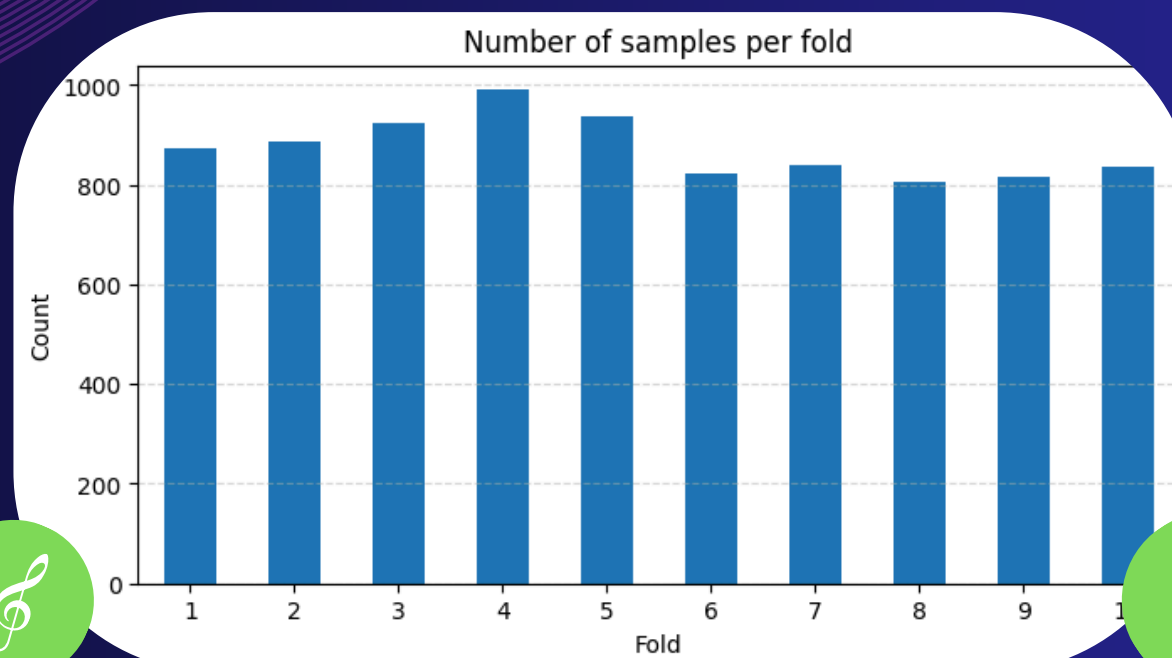
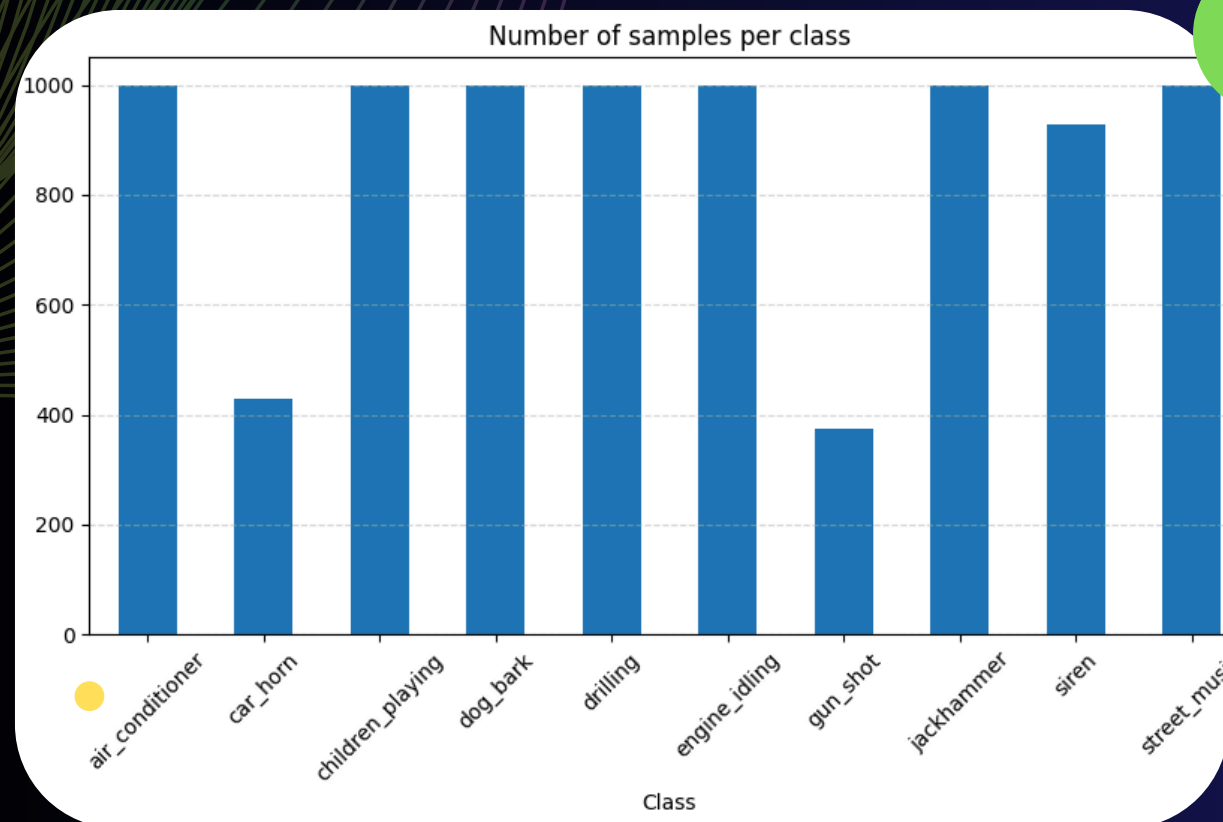
Aprendizagem Computacional II – LIACD 2025/26

Autores: Alexandre Furriel, Daniel Gomes e Liliana Silva



Características do Dataset

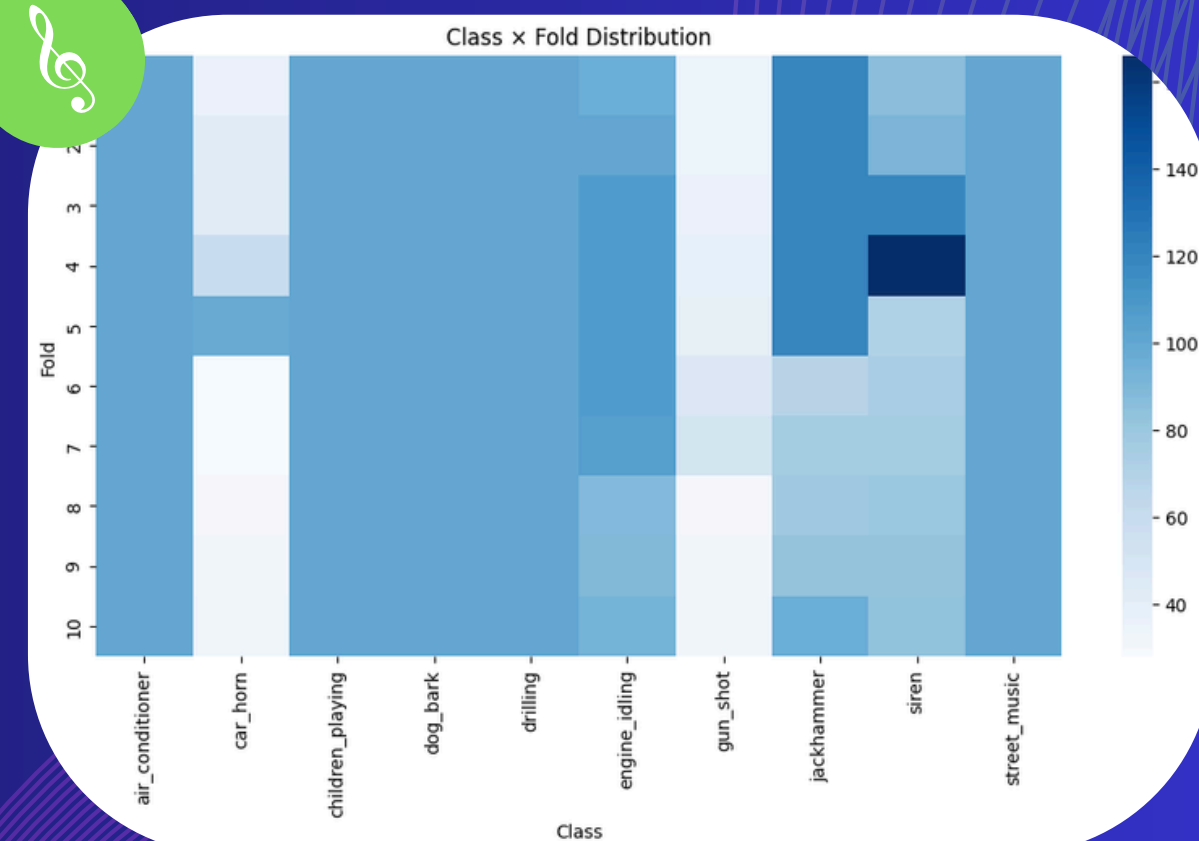
- 8732 excertos de áudio ($\leq 4s$)
- 10 classes de sons urbanos
- Divisão oficial em 10 folds
- Variabilidade elevada entre samples (ambiente, volume, ruído)



 Dataset

Padrões observados na análise

- Desbalanceamento claro em classes como car_horn e gun_shot
- Tamanho dos folds moderadamente desigual
- Distribuição classe \times fold bastante estável



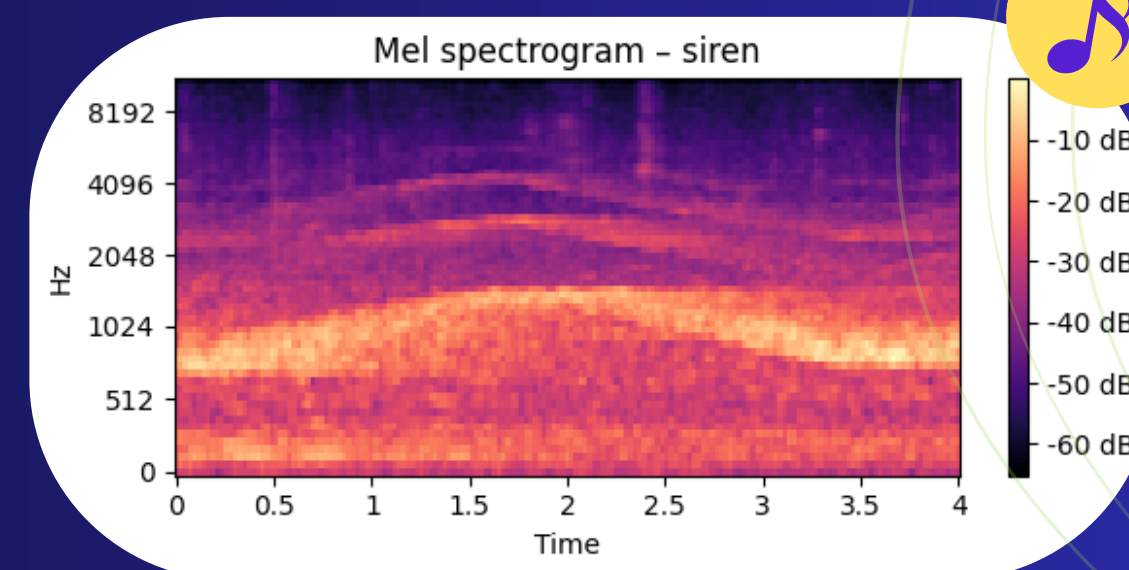
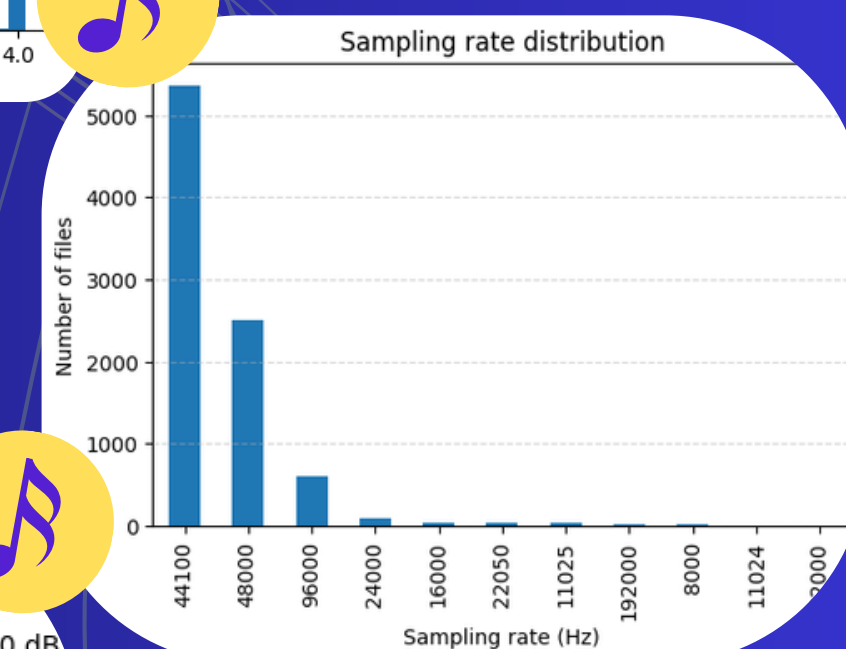
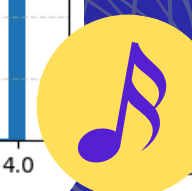
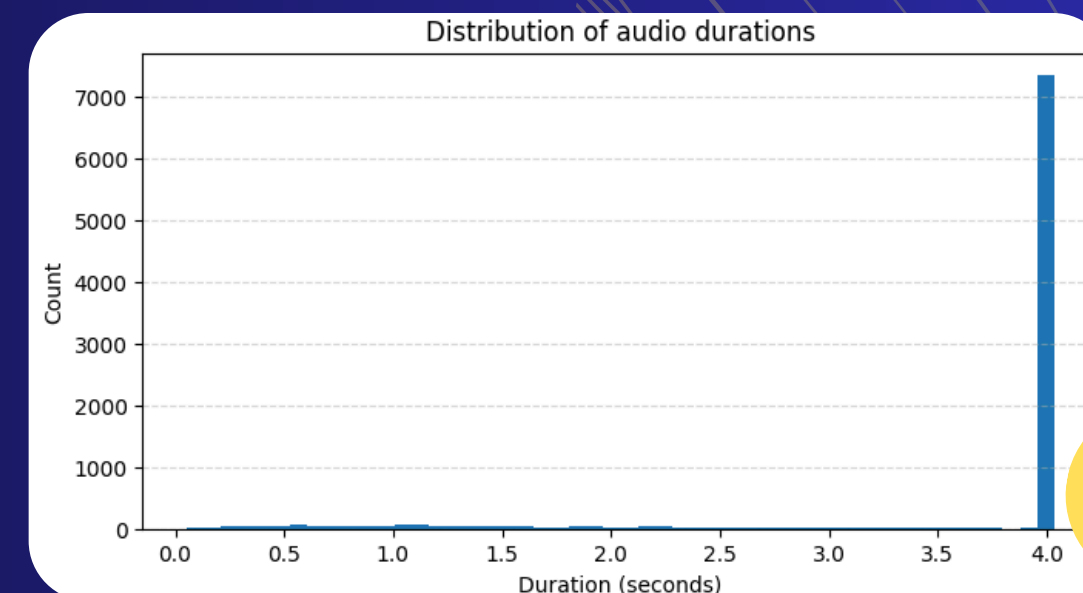
🎵 Pré-Processamento

Desafios

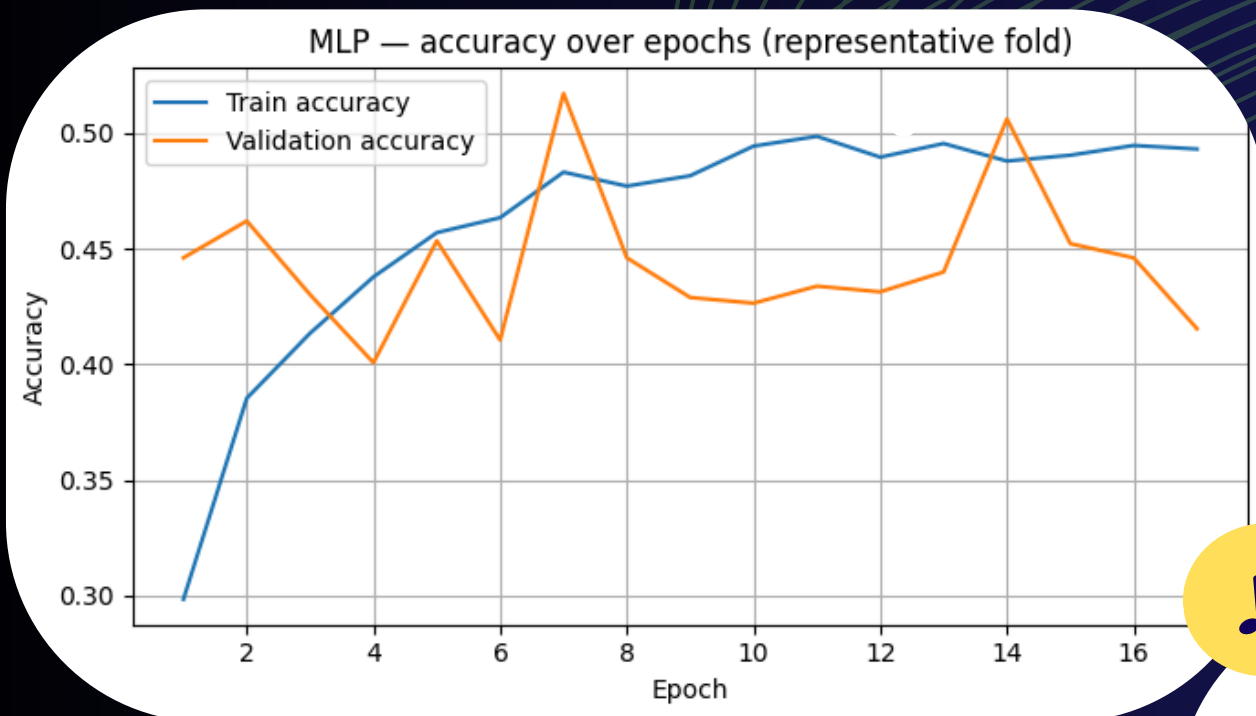
- Durações inconsistentes
- Sampling rates diferentes
- Variabilidade de ruído e volume

Decisões

- Uniformização da sampling rate
- Padding para duração fixa
- Normalização dos sinais
- Extração de Mel Spectrograms
- Preparação de tensores para MLP e CNN

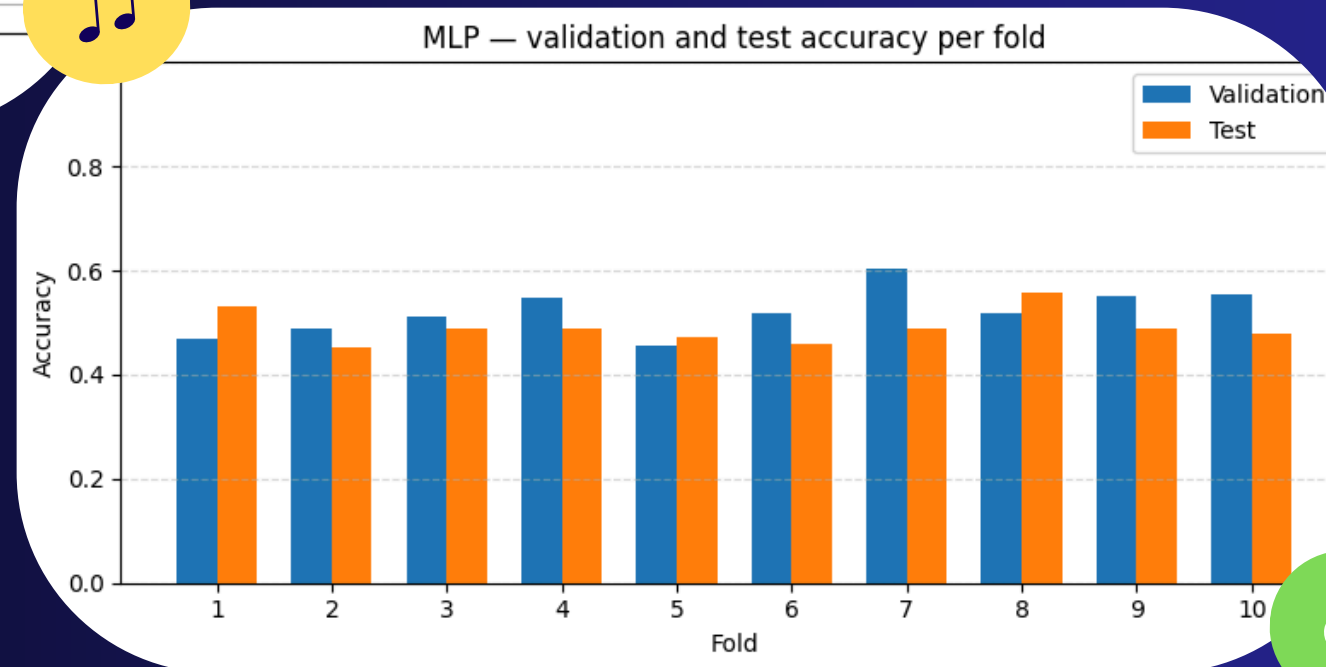


MLP — Modelo Baseline



Input

- MelMel spectrogram (64×173) flatten \rightarrow vetor 1D ($\sim 11k$ features)



Comportamento

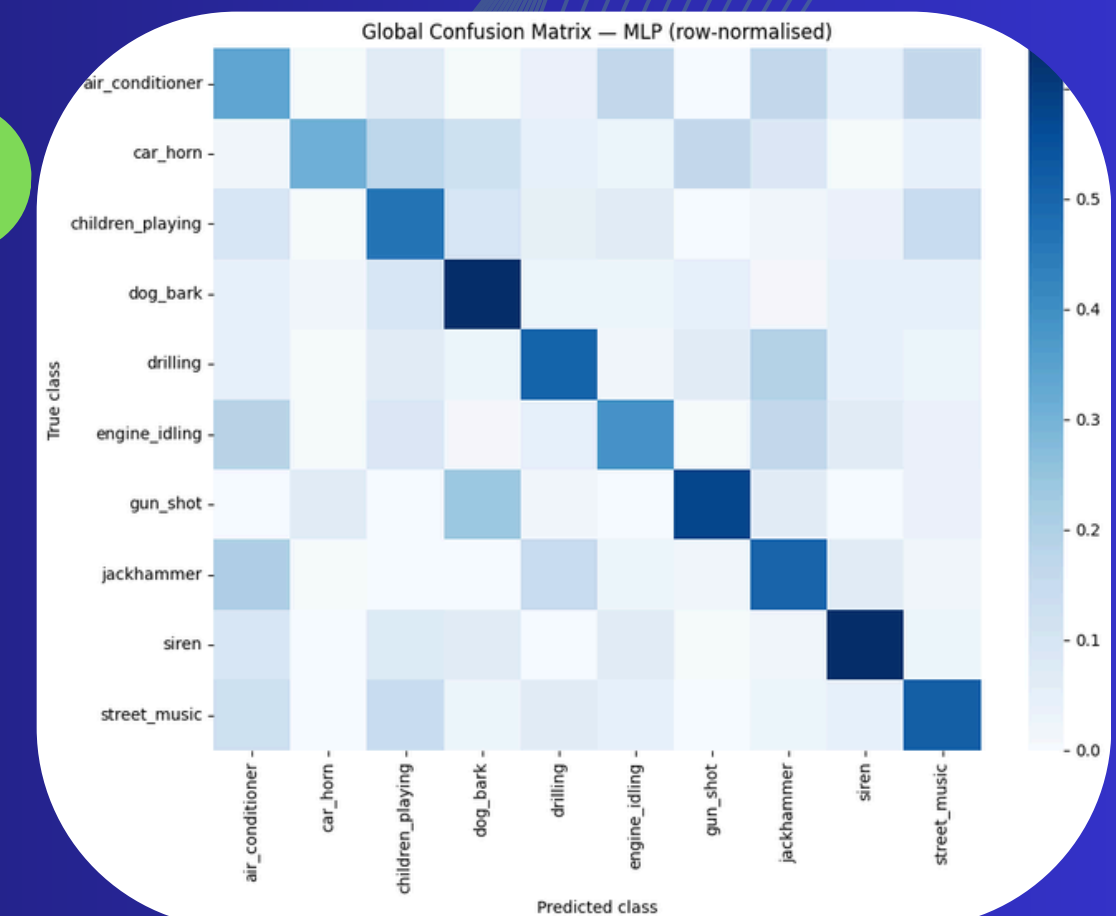
- Convergência rápida
- Overfitting moderado (train > val)
- Validação instável (flutuações elevadas)

Desempenho (10-fold)

- Val acc $\approx 53\% \pm 5\%$
- Test acc $\approx 48\% \pm 5\%$
- Erros sistemáticos em classes contínuas (air_conditioner, engine_idling)

Arquitetura

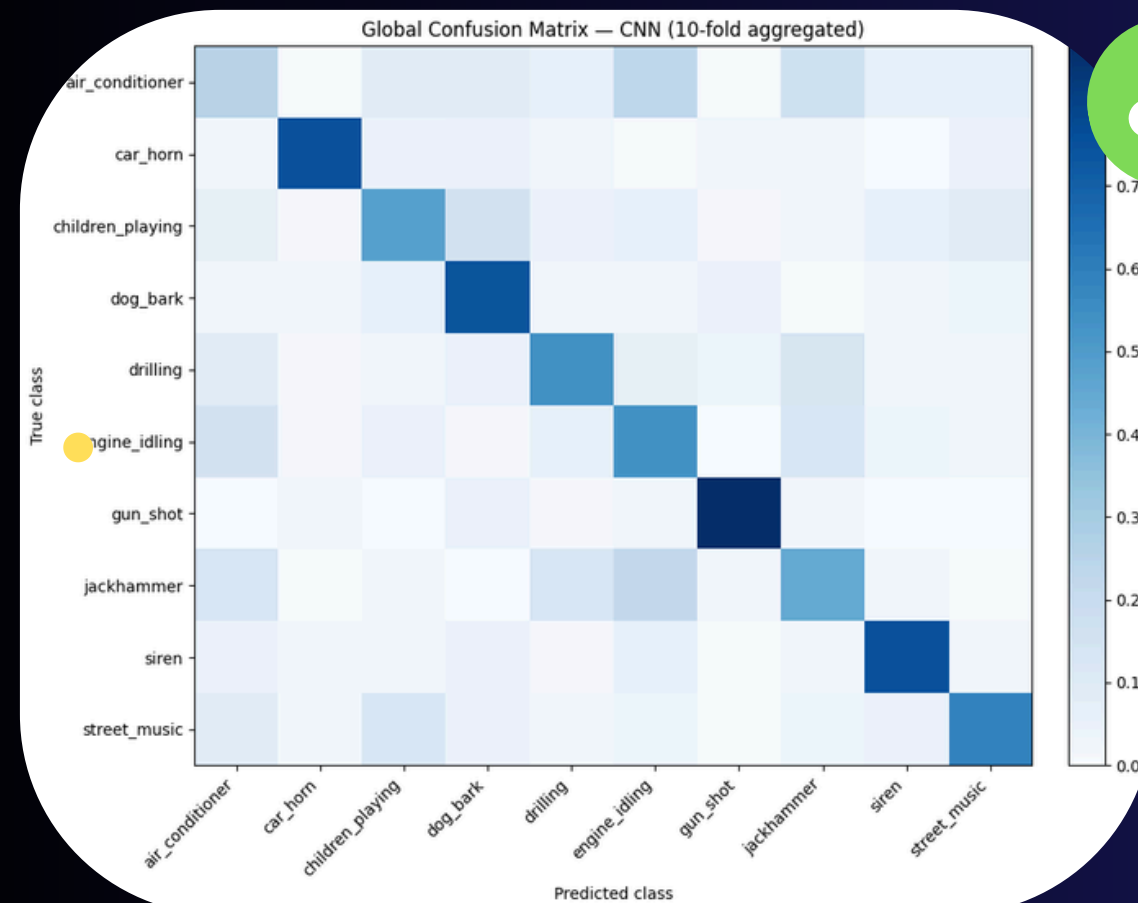
- 2 camadas Dense ($512 \rightarrow 256$) + ReLU
- Batch Normalization
- Dropout (0.5)
- Softmax para 10 classes



CNN (64 Mel Bands)

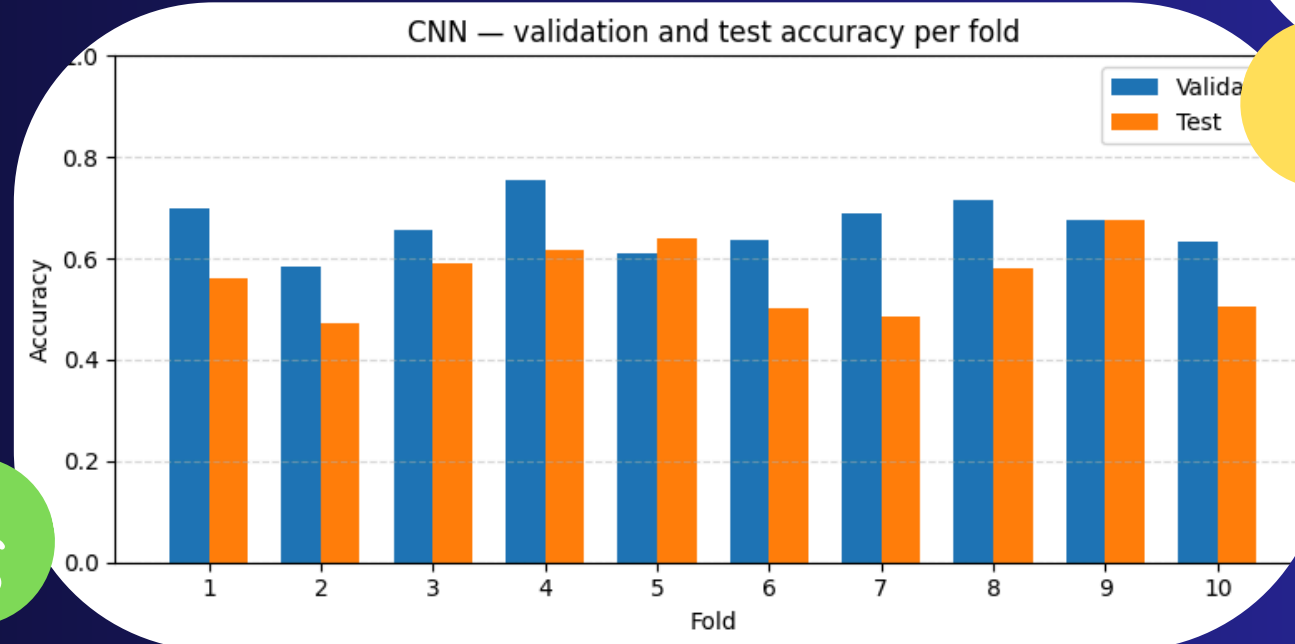
Arquitetura

- Conv2D + BatchNorm + ReLU
- MaxPooling
- Dropout
- Dense final (10 classes)



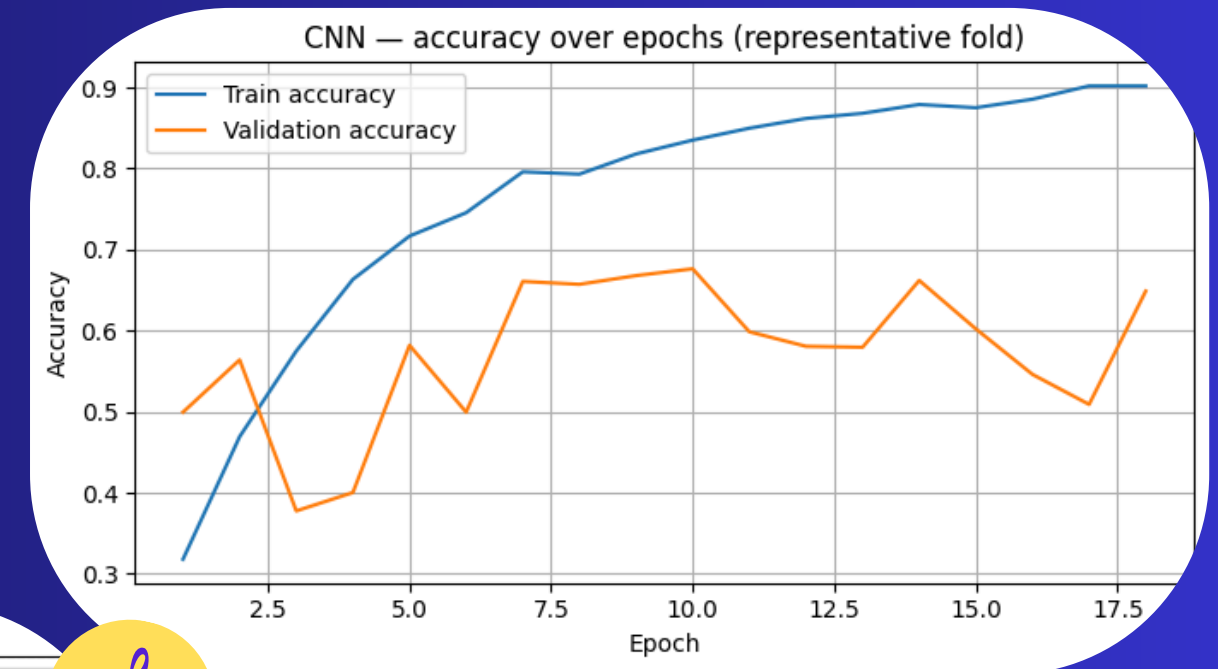
Input

- Mel spectrogram 2D ($64 \times 173 \times 1$)



Desempenho (10-fold)

- Val acc $\approx 0.65 \pm 0.05$
- Test acc $\approx 0.58 \pm 0.06$
- Grandes melhorias em classes transitórias (car_horn, dog_bark)
- Confusão ainda presente entre classes contínuas (air_conditioner, engine_idling)



Comportamento

- Melhor generalização vs MLP
- Menor overfitting
- Validação estável mas com flutuações iniciais
- Captura padrões locais (transientes e texturas)



Refinamento Espectral: 64 → 96 Mel Bands

Porque ajustar:

- Baixa resolução nos 64 bins limitava a separação em regiões críticas de baixa frequência.
- Harmônicos mecânicos ficavam sobrepostos e difíceis de distinguir.
- Cenas amplas exigiam maior granularidade espectral.

Decisão:

- Aumentar a resolução para 96 Mel bands, mantendo todo o resto constante (modelo, treino, folds).

Objetivo:

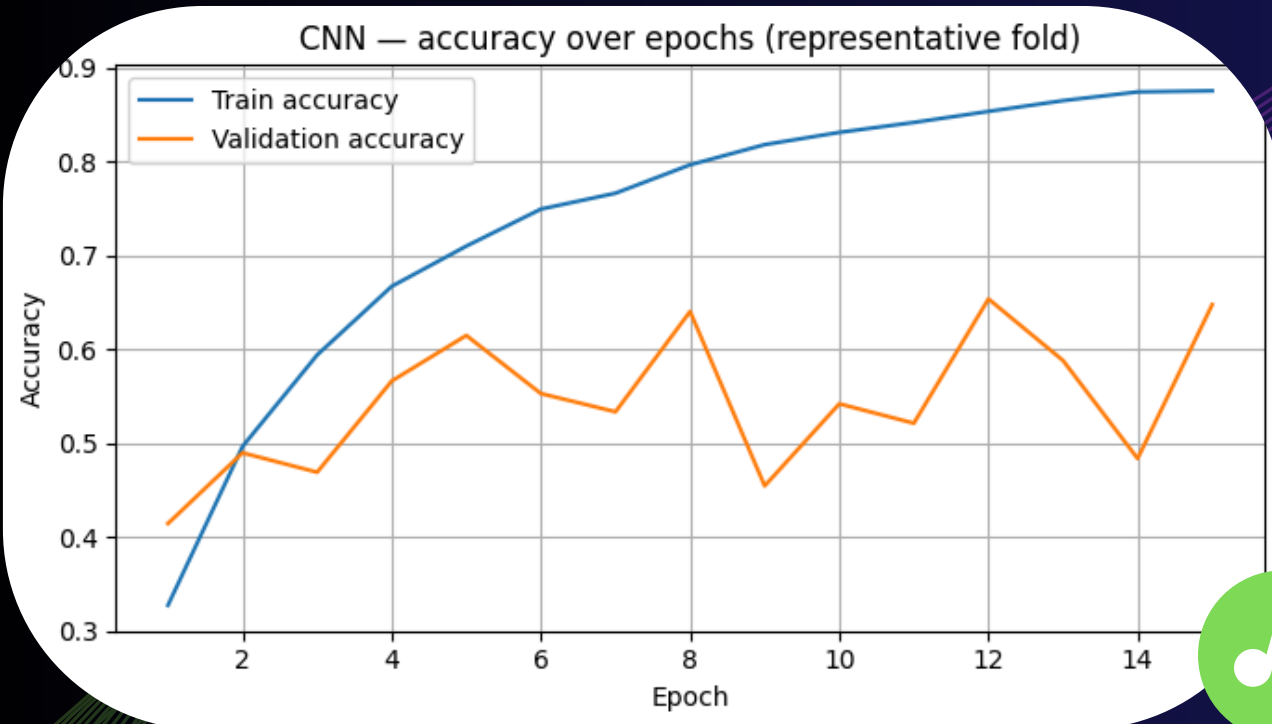
- Verificar se o aumento de detalhe espectral reduz a confusão e melhora a generalização.

Problema Detetado no CNN-64

| | |
|---------------------------|----------------------------------------------------------------------------------------------------|
| Baixa frequência | Confusão entre <i>air_conditioner</i> e <i>engine_idling</i> devido a espectros quase idênticos |
| Sons mecânicos periódicos | <i>drilling</i> ↔ <i>jackhammer</i> apresentam harmônicos sobrepostos |
| Cenas complexas | <i>children_playing</i> ↔ <i>street_music</i> apresentam mistura de fontes e variabilidade elevada |



CNN (96 Mel Bands)



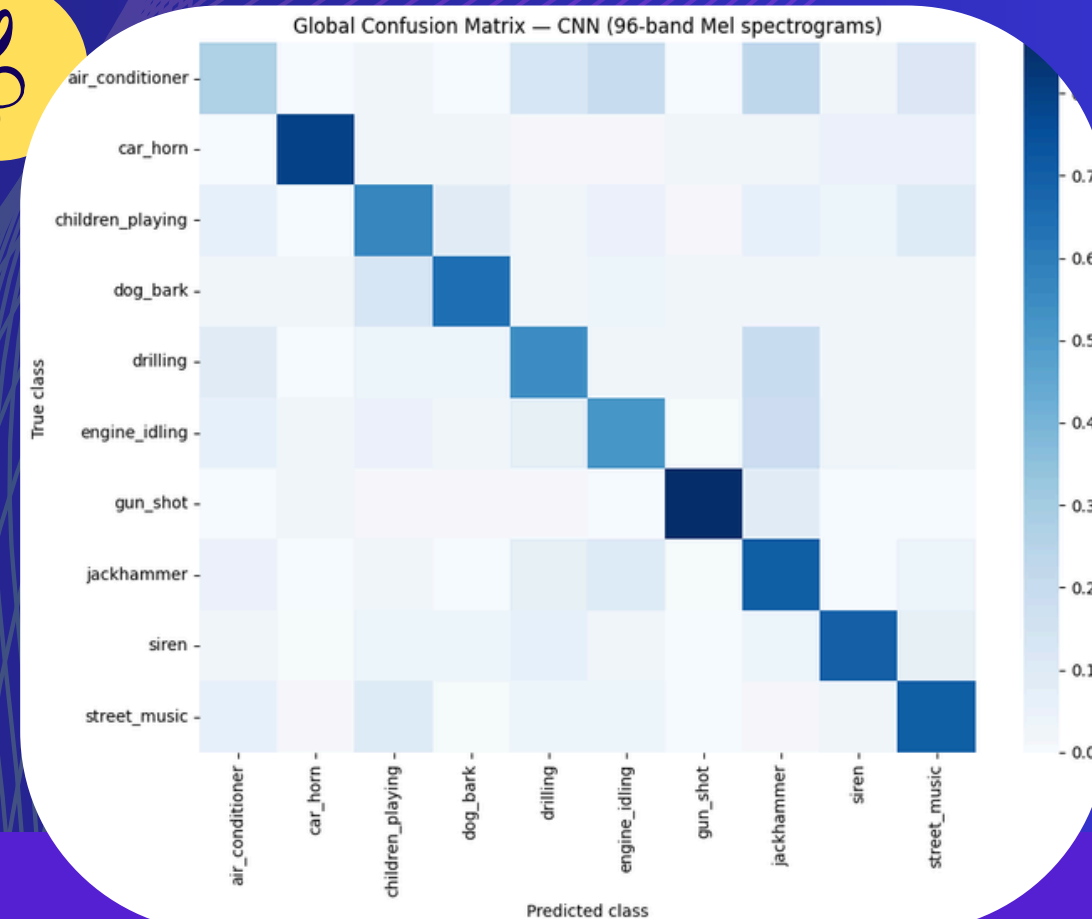
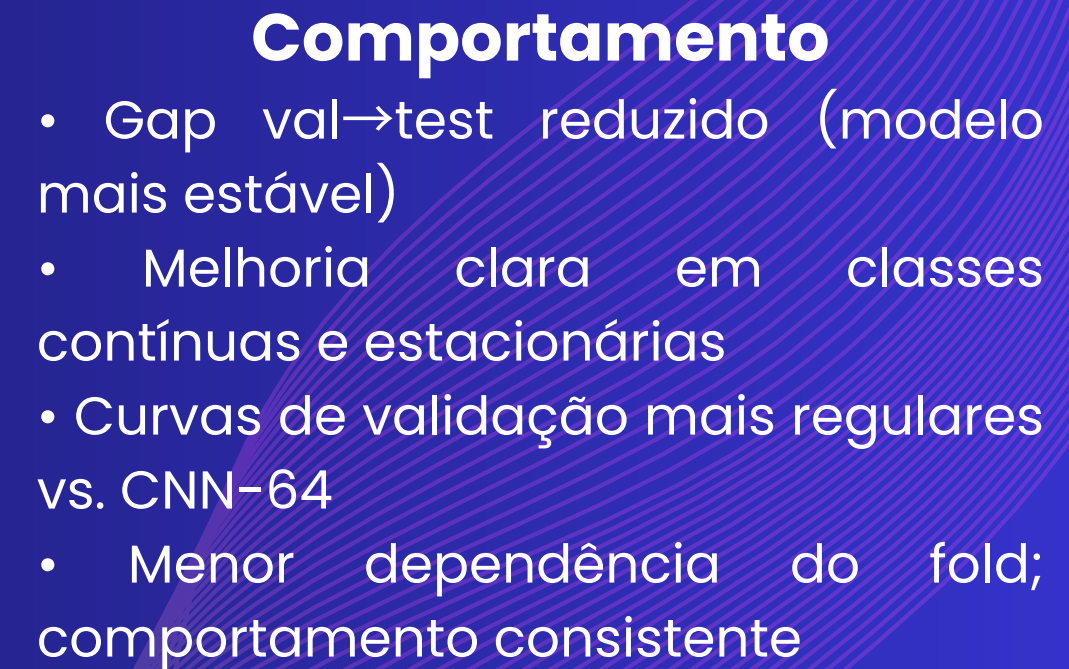
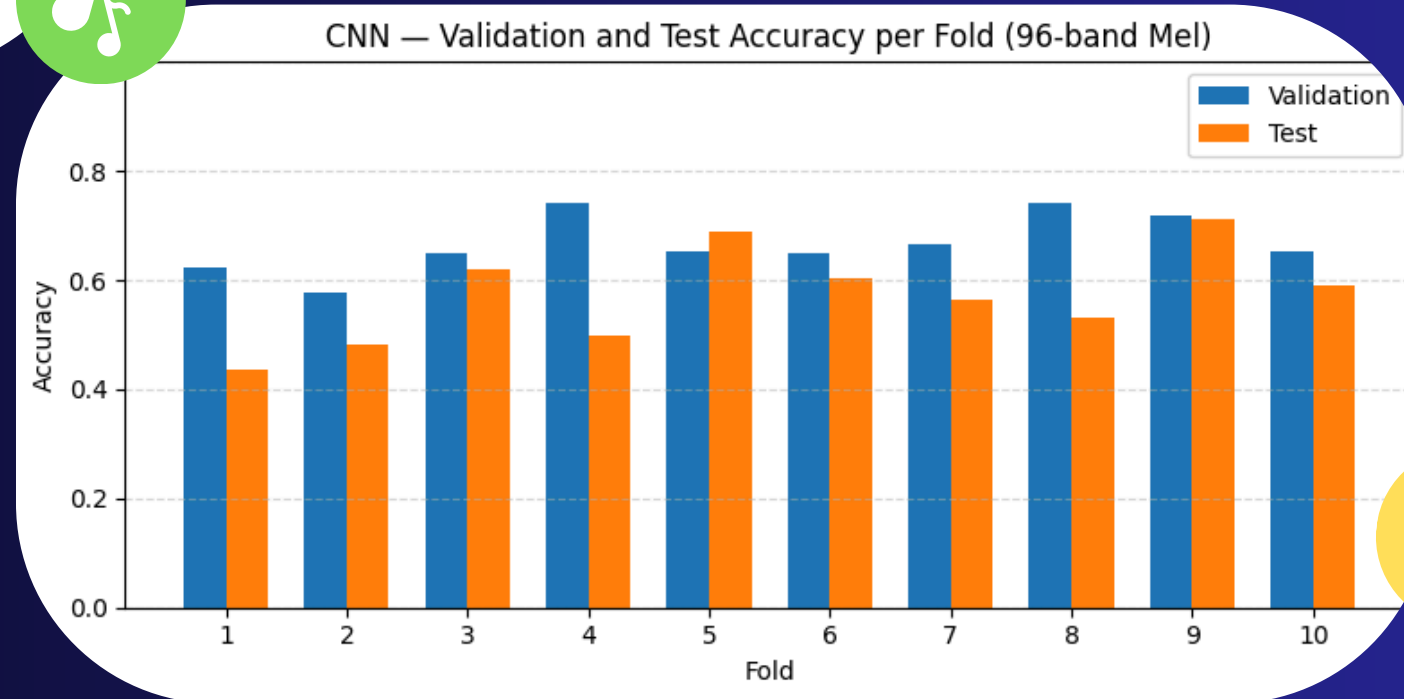
Desempenho (10-fold)

- Val acc $\approx 0.687 \pm 0.045$
- Test acc $\approx 0.607 \pm 0.049$
- Variância do teste ≈ 0.049
- Ganhos particularmente fortes em classes mecânicas

Insights da Matriz de Confusão

- Diagonais mais fortes em classes de baixa frequência
- Melhor separação em padrões periódicos (drilling, jackhammer)
- Menor leakage children_playing ↔ street_music
- Confusão persistente apenas em pares acusticamente semelhantes

- Mel spectrogram 2D ($96 \times 173 \times 1$)
- Mesma arquitetura do CNN-64





CNN-96 + Reg/Aug

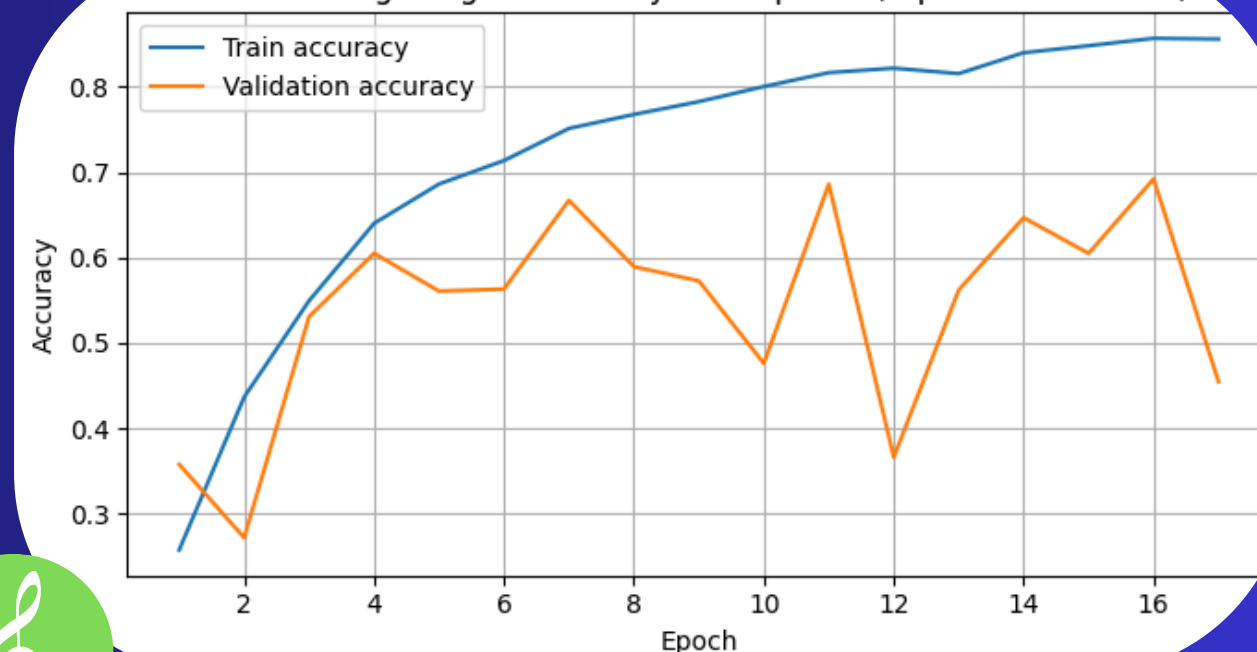
Comportamento

- Diferença val→test mais pequena → overfitting reduzido
- Curvas de treino/validação mais estáveis
- Maior consistência entre folds → modelo mais robusto

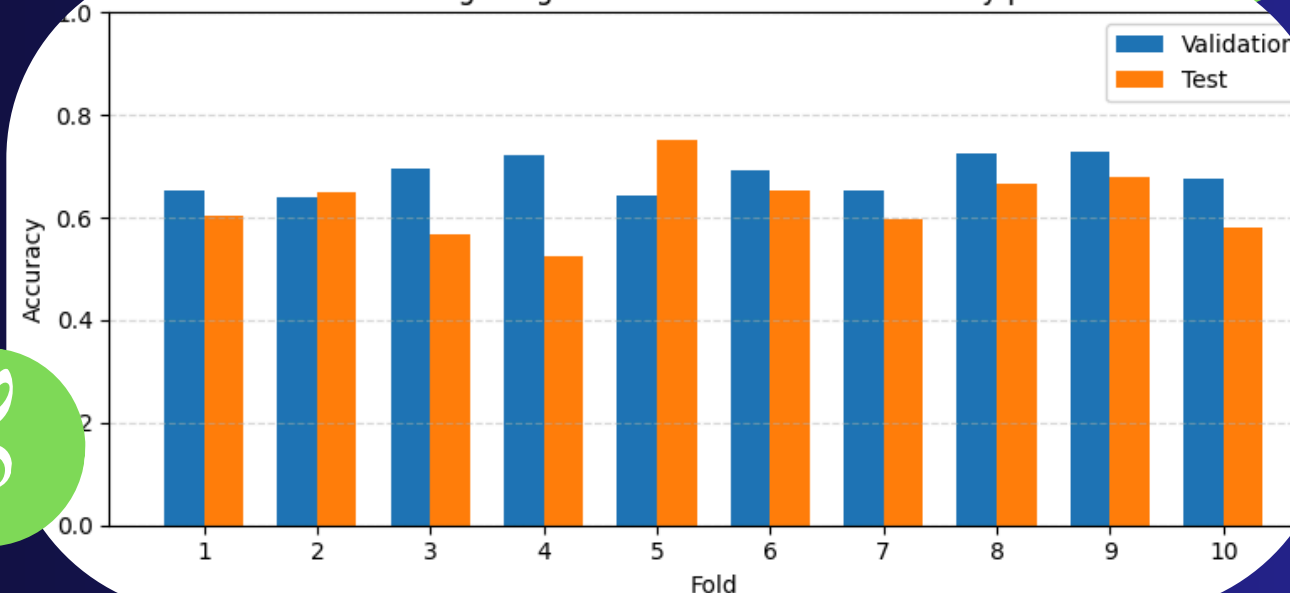
Refinamentos aplicados

- Dropout & L2 mais fortes
- Light SpecAugment (time + freq masking) apenas no treino
- Arquitetura idêntica → alteração apenas na variabilidade do treino

CNN-96 Reg+Aug — accuracy over epochs (representative fold)



CNN-96 Reg+Aug — Validation and Test Accuracy per Fold

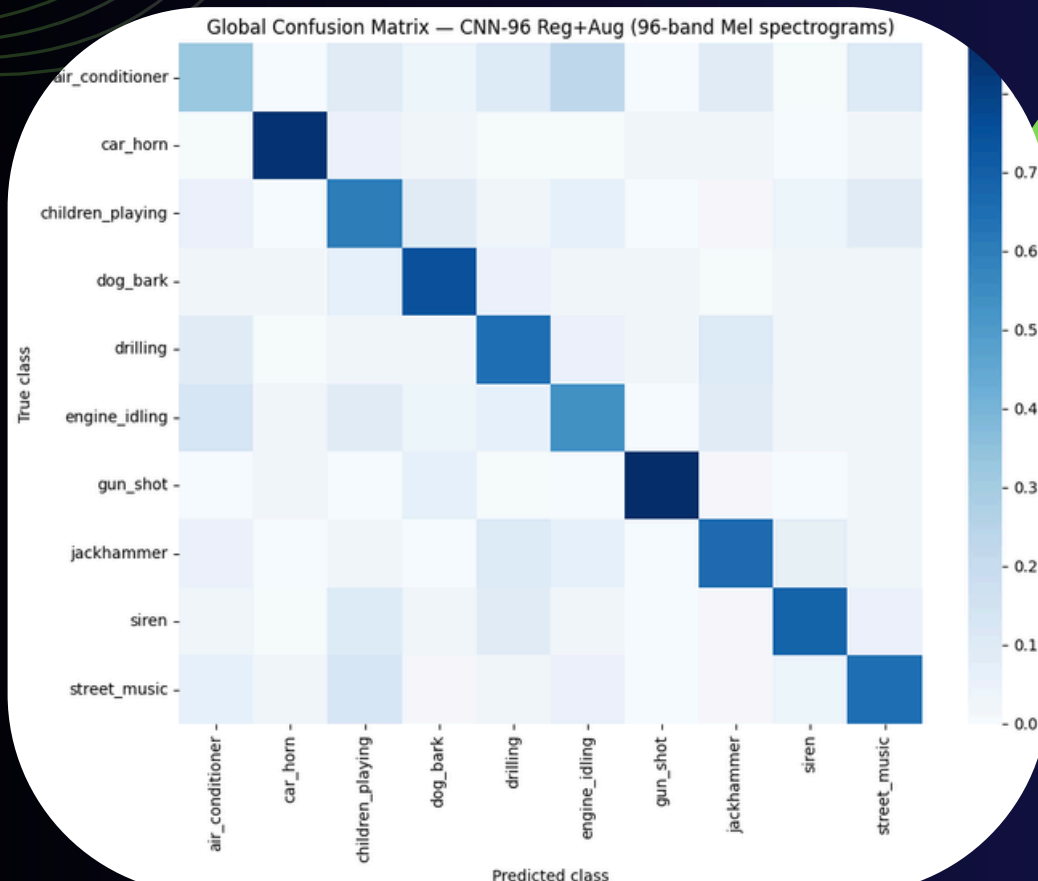


Desempenho (10-fold)

- Val acc $\approx 0.679 \pm 0.028$
- Test acc $\approx 0.620 \pm 0.047$
- Variância do teste mais baixa (≈ 0.047)
- Pequeno trade-off: ligeira queda em validação, mas melhor generalização

Insights da Matriz de Confusão

- Melhoria adicional em classes mecânicas (drilling, jackhammer)
- Menor leakage residual em children_playing ↔ street_music
- Confusão quase apenas nos pares acusticamente semelhantes



Comparação Final dos Modelos

| Modelo | Val acc (média \pm std) | Test acc (média \pm std) | Var. do teste | Características principais |
|------------------|------------------------------|-------------------------------|------------------|------------------------------------------------------------------------------------|
| MLP | 0,5304 \pm 0,0511 | 0,4778 \pm 0,0529 | \sim 0,053 | Sem estrutura 2D; forte overfitting; má separação entre classes. |
| CNN 64 bandas | 0,6688 \pm 0,0385 | 0,5387 \pm 0,1069 | \sim 0,107 | Explora o espectrograma 2D; boa em eventos curtos; variância alta entre folds. |
| CNN 96 bandas | 0,6872 \pm 0,0452 | 0,6072 \pm 0,0490 | \sim 0,049 | Melhor resolução em baixas frequências; ganhos em classes contínuas; mais estável. |
| CNN 96 + Reg/Aug | 0,6788 \pm 0,0276 | 0,6197 \pm 0,0466 | \sim 0,047 | Regularização + masking; melhor generalização global; menor sensibilidade ao fold. |

- O aumento da resolução Mel (64 \rightarrow 96) foi a melhoria mais significativa, resolvendo o principal bottleneck representacional.
- A regularização + masking reduziu o overfitting residual e estabilizou os resultados sem modificar a arquitetura.
- O modelo final (CNN-96 + Reg/Aug) apresenta o melhor equilíbrio entre desempenho, generalização e robustez.

• Bónus e Desafios

Bónus – DeepFool

- Aplicámos o ataque DeepFool à CNN-96 Reg+Aug no último fold de validação.
- Avaliámos 100 amostras de teste.
- O ataque enganou o modelo em 87% dos casos.
- A perturbação média tinha norma $L_2 \approx 1.38 \times 10^6$, praticamente impercetível no Mel-spectrograma.
- Mostra que, apesar da boa accuracy, o modelo mantém baixa robustez adversarial.

Desafios Identificados

- O modelo é sensível a pequenas perturbações → baixa robustez adversarial.
- Classes com padrões espectrais semelhantes continuam difíceis (ex.: air_conditioner vs engine_idling).
- Variância entre folds mostra dependência do conjunto de treino.
- Limitações computacionais impediram explorar arquiteturas maiores ou augmentations avançadas.