

# Transfer Learning for Smart Background Simulation at Belle II

<u>David Giesegh</u>, Nikolai Krug, Boyang Yu, Thomas Kuhr LMU Munich

> 12.03.2025 KISS annual meeting









# nulation

#### Smart Background Simulation

(Introduced in PhD Thesis by James Kahn (2019))



- Skims are centrally produced loose selections tailored to specific analysis types
   → Many skims throw away large fractions of costly simulated background events
- Event generation much faster than detector simulation and event reconstruction  $(t_{slow}/t_{fast} \sim 100)$
- Typically, no simple correlations between features after generation and final selection criteria
- Perhaps a **neural network based classifier** can solve this issue?

#### Filtering procedure



- False positives: unproblematic (simply lower performance)
- False negatives: problematic! Discarded events are gone and may thus bias distributions
- Solution: importance sampling



#### Importance sampling

(Introduced in Master's Thesis by Boyang Yu (2021))

Idea:

- use NN output as sampling probability p<sub>i</sub>
- weight each event by its inverse probability  $\omega_i$
- > No bias by construction (Disadvantage: weighted events with arbitrarily high weights)

Propapility	
Event 1	0.75
Event 2	0.33
Event 3	0.01
Event 4	0.06
Event 5	0.50

![](_page_3_Picture_8.jpeg)

![](_page_3_Picture_9.jpeg)

Event 1	1.33
Event 2	3.00
Event 3	X
Event 4	16.67
Event 5	X

Weight

#### Evaluation metric: speedup

- Goal: Train NN to achieve highest **speedup**  $(t_{noNN}/t_{NN})$  for the same effective sample size
- Effective sample size:
  - := size of an unweighted sample that would yield the same relative statistical uncertainty
  - Estimate for number of expected events:  $N_{exp} = \sum_i \omega_i$
  - Estimate for variance:  $\sigma^2 = \sum_i \omega_i^2$ • Effective sample size:  $\frac{\sqrt{N_{eff}}}{N_{eff}} = \frac{\sqrt{\sum_i \omega_i^2}}{\sum_i \omega_i} \Rightarrow N_{eff} = \frac{(\sum_i \omega_i)^2}{\sum_i \omega_i^2}$

#### Evaluation metric: speedup

- Goal: Train NN to achieve highest **speedup**  $(t_{noNN}/t_{NN})$  for the same effective sample size
- Speedup also depends on:
  - Assumed times for generation  $t_{fast}$ , NN inference  $t_{NN}$  and simulation/reconstruction  $t_{slow}$  (roughly  $t_{fast}$ :  $t_{NN}$ :  $t_{slow} \sim 1 : 1 : 100$ )
  - Skim efficiency (retention rate)
    - $\rightarrow$  lower retention increases potential speedup

![](_page_6_Picture_0.jpeg)

### Evaluation metric: speedup

- Speedup =  $t_{noNN}/t_{NN}$
- Define:
  - $t_1 \coloneqq t_{fast} + t_{NN} + t_{slow}$  (events accepted by NN)
  - $t_2 \coloneqq t_{fast} + t_{NN}$  (events rejected by NN)
  - $t_3 \coloneqq t_{fast} + t_{slow}$  (events with no NN)
- Use:
  - $t_{noNN} = t_3 \cdot N_{eff}$
  - $t_{NN} = t_1 \cdot (N_{TP} + N_{FP}) + t_2 \cdot (N_{FN} + N_{TN})$
- Some algebra later...

	Pass skim	Fail skim
NN accepts	ТР	FP
NN rejects	FN	TN

Speedup = 
$$\frac{t_3/E[\omega_i]}{(\varepsilon f_{TP} + (1 - \varepsilon) f_{FP}) t_1 + (\varepsilon f_{FN} + (1 - \varepsilon) f_{TN}) t_2}$$

![](_page_7_Figure_0.jpeg)

#### Optimizing for speedup

![](_page_7_Figure_2.jpeg)

- Model trained with **cross entropy loss** and final sigmoid activation
- Idea: optimal classification probability for cross entropy should be related to optimal sampling probability for speedup by monotonic function
- $\rightarrow$  Replace sigmoid activation by fitted "skewed" sigmoid optimized for speedup

#### Optimizing for speedup

- Turns out generic skewed sigmoid  $(1 + ae^{-c(x-b)})^{-1/a}$  has too many parameters (no unique maximum)
- Tradeoff: shorter time  $\leftrightarrow$  narrower weight distribution
- Narrow weight distribution desirable → empirically leads to the following function:

![](_page_8_Figure_4.jpeg)

#### Dataset

- Generator level MC particles (hadrons, no gluons/quarks)
- List of particles with decay structure information (mother-daughter relations)
- Particle features:
  - PID
  - 4-momentum
  - Production 4-vertex

![](_page_9_Figure_7.jpeg)

#### Model architecture

- Based on ParT (Particle Transformer for Jet Tagging, arxiv:2202.03772)
   → State-of-the-art performance in jet tagging by pre-training on large dataset and subsequent fine-tuning (e.g. for top tagging)
- Very generic architecture (essentially just a transformer)
- Supports edge features:
  - Adjacency matrix of decay graph
  - Two-particle invariant masses and angles
- 10 layers, ~2 million parameters

![](_page_11_Figure_0.jpeg)

![](_page_11_Picture_1.jpeg)

![](_page_12_Picture_0.jpeg)

#### Transfer learning

- Fundamental problem:
  - Training transformer requires large dataset
  - But if a large dataset for a skim is already available: what do I need the model for?
- Solution: pre-trained models!
- Two approaches:

#### Feature extraction:

- Remove final layer and retrain only that
- Perhaps even just re-fit final sigmoid activation
- Only feasible for highly correlated skims

#### Whole model fine-tuning:

- Start with pre-trained model but adjust all parameters
- Perhaps reinitialize last layer (if output shape changes)
- Promising results in ParT paper

#### Transfer learning - Prestudy

![](_page_13_Figure_1.jpeg)

- Model trained on reference skim, then fine-tuned on 2000 events of new skim
- Feature extraction works poorly, especially with low correlations between skims (A and B)
- Training from scratch overtrains very quickly on small dataset
- Full fine-tuning looks promising, consistently yielding the best validation loss

#### Large scale training

- Idea: to make flexible fine-tuning possible, pre-train model on all available skims at once
  - Currently samples for 51 centrally produced skims available
- 7 background types (useful for many analyses):  $e^+e^- \rightarrow B^+B^-$ ,  $B^0\overline{B}^0$ ,  $q\overline{q}$  (q=u,d,s,c),  $\tau^+\tau^-$ 
  - Condition model on background type
- Dataset balanced between different types, containing  ${\sim}180$  million events or  ${\sim}20 fb^{-1}$
- Training
  - Without class weights (original)
  - With class weights
  - With class weights and empirical *focal loss* emphasizing difficult training examples

#### Training curves

![](_page_15_Figure_1.jpeg)

#### Speedups

![](_page_16_Figure_1.jpeg)

![](_page_16_Picture_2.jpeg)

Achievable speedups correlate with

• Separation power (as measured by the area under the ROC curve)

#### • Skim efficiency

(lower retention rates  $\rightarrow$  higher potential for speedup)

#### 1e7 1.4 Intelligent Standard 1.2 --- N<sub>eff</sub> Time Consumption in ms 1.0 0.8 0.6 0.4 0.2 0.0 1000 2000 3000 0 4000 No. of Events Passed through Entire Chain

## Adaptive/Reinforcement Learning

(Introduced in Bachelor's Thesis by Daniel Pollmann (2024))

- Another possible idea:
  - Train model while producing data and running skim
  - Model becomes successively better, producing data more efficiently
- Pro: Overall time-saving, on-the-fly procedure in one step
- Contra: Requires implementing training loop in Belle II production software

#### Current work and outlook

![](_page_18_Picture_1.jpeg)

Currently three main objectives:

- Test pre-trained model on already known skim under different data taking conditions
  - E.g. reprocessed data under a new software release or with more realistic beam background from so-called run-dependent MC
- Fine-tune and evaluate pre-trained model on entirely new skim
- Implement model inference in Belle II software to test model in an actual skim production and evaluate realistic speedup

![](_page_19_Picture_0.jpeg)

![](_page_19_Picture_1.jpeg)

# Thank you for your attention!

KISS annual meeting - 12.03.2025

![](_page_20_Picture_0.jpeg)

# Backup

KISS annual meeting - 12.03.2025

#### Speedup derivation

• 
$$t_{noNN} = t_3 \cdot N = t_3 \cdot \frac{(\sum_i \omega_i)^2}{\sum_i \omega_i^2} = t_3 \cdot \frac{(\varepsilon N_{gen})^2}{\varepsilon^2 N_{gen} E[\omega_i]} = \frac{t_3 N_{gen}}{E[\omega_i]}$$
  
•  $t_{NN} = t_1 \cdot (N_{TP} + N_{FP}) + t_2 \cdot (N_{FN} + N_{TN}) = t_1 \cdot N_{gen} (\varepsilon f_{TP} + (1 - \varepsilon) f_{FP}) + t_2 \cdot N_{gen} (\varepsilon f_{FN} + (1 - \varepsilon) f_{TN})$ 

Speedup = 
$$\frac{t_3/E[\omega_i]}{(\varepsilon f_{TP} + (1 - \varepsilon) f_{FP}) t_1 + (\varepsilon f_{FN} + (1 - \varepsilon) f_{TN}) t_2}$$

![](_page_21_Picture_3.jpeg)

#### Hyperparameters

- Largely based on ParT:
  - 8 self attention blocks
  - 2 class attention blocks
  - 8 heads per multi-head attention block
  - Embedding size 128
  - $\rightarrow$  Roughly 2 million parameters
- Differences:
  - Fewer norm layers
  - Embeddings for PDG ID and background type
  - Using 3 pair features (adjacency matrix, pair mass and angle)

![](_page_22_Picture_11.jpeg)