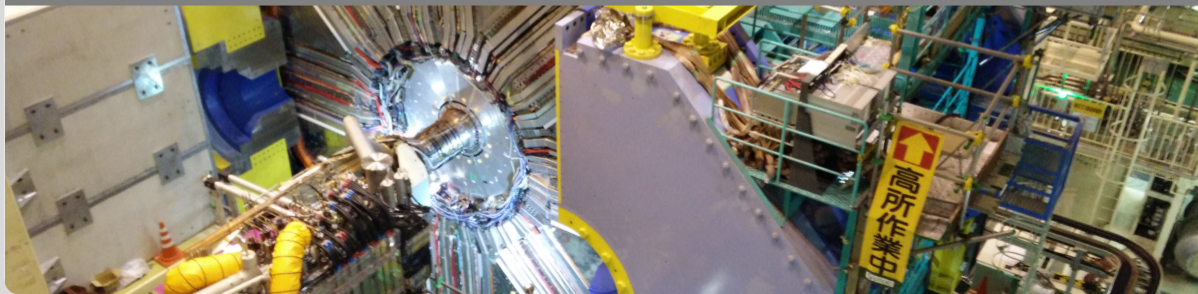


Selective background Monte Carlo simulation at Belle II

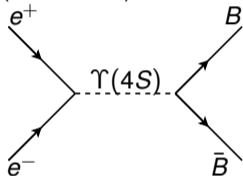
James Kahn, Andreas Lindner, Thomas Kuhr | 5th November 2019

INSTITUT FÜR EXPERIMENTELLE TEILCHENPHYSIK (ETP)

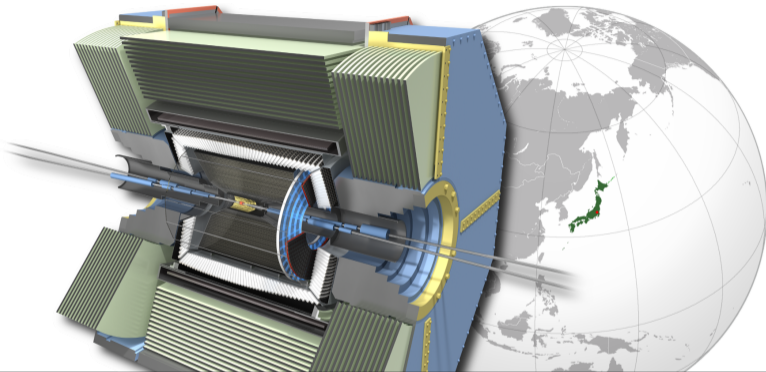


Belle II Experiment

Asymmetric e^+e^-
experiment mainly at
the $\Upsilon(4S)$ resonance
(10.58 GeV)



Focus on B, charm
and τ physics



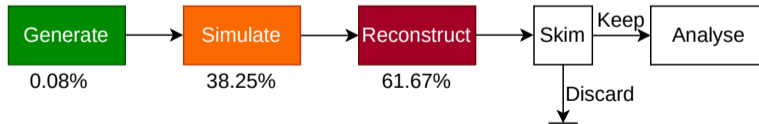
	KEKB/Belle	SuperKEKB/Belle II
Operation	1999–2010	2019–2027
Peak luminosity	$2.11 \times 10^{34} \text{ cm}^{-2}\text{s}^{-1}$	$8 \times 10^{35} \text{ cm}^{-2}\text{s}^{-1}$
Integrated luminosity	1 ab^{-1} (772 million $B\bar{B}$ pairs)	50 ab^{-1}

Problem

- Approach at Belle:
 - Background MC $\approx 10 \times$ data
- Infeasible at Belle II \rightarrow still require high statistics
- Currently: ~ 100 HS06 s/event
 - $1 \text{ ab}^{-1} \approx 80 \text{ GHS06 s}$

Skims

- Physics working-group specific datasets (26)
- General selections applied to discard trivial backgrounds
- Retain $\mathcal{O}(0.1\text{--}10\%)$ of full dataset



Problem

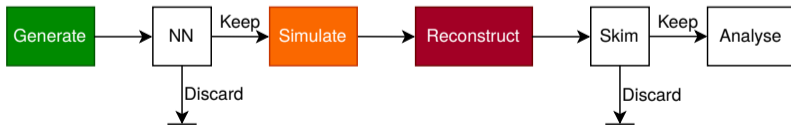
- Approach at Belle:
 - Background MC $\approx 10 \times$ data
- Infeasible at Belle II \rightarrow still require high statistics
- Currently: ~ 100 HS06 s/event
 - $1 \text{ ab}^{-1} \approx 80 \text{ GHS06 s}$

Proposed solution:

Insert NN to predict skims before expensive steps

Skims

- Physics working-group specific datasets (26)
- General selections applied to discard trivial backgrounds
- Retain $\mathcal{O}(0.1\text{--}10\%)$ of full dataset



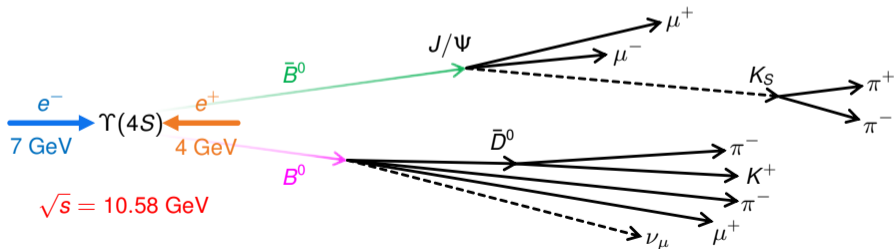
Dataset

~ 300,000 particle collision events with binary classification labels:

- Hadronic B^+ meson reconstruction ($\sim 5\%$)
- Time-dependent CP violation ($\sim 0.2\%$)

Graph terminology

- Nodes = Particles
- Node attributes = Particle properties
- Edges = Parent-daughter relations (decays)
- Graph type = Tree



~ 300,000 particle collision events with binary classification labels:

- Hadronic B^+ meson reconstruction (~ 5%)
- Time-dependent CP violation (~ 0.2%)

Graph terminology

- Nodes = Particles
- Node attributes = Particle properties
- Edges = Parent-daughter relations (decays)
- Graph type = Tree

$\Upsilon(4S)$ (300553)
 B^0 (-511)
 J/ψ (443)
 μ^+ (-13)
 μ^- (13)
 K_S^0 (310)
 π^- (-211)
 π^+ (211)
 B^0 (511)
 \bar{D}^0 (-421)
 π^- (-211)
 K^+ (321)
 π^- (-211)
 μ^+ (-13)
 ν_μ (14)

Feature	Definition
PDG code	Identifier of particle type and charge.
Mother PDG code	Particle parent PDG code.
Mass	Particle mass in GeV/c^2 .
Charge	Electric charge of the particle.
Energy	Particle energy in GeV.
Momentum	Three momentum of the particle in Gev/c .
Production time	Production time in ns relative to $\Upsilon(4S)$ production.
Production vertex	Coordinates of particle production vertex.
Status bit	Bitmask representing MC production conditions.

Graph Isomorphism Network

Node N update rule of layer ℓ (Red = trainable):

$$N^{(\ell+1)} = \text{MLP}^{(\ell)} \left(W_p^{(\ell)} N_p^{(\ell)} + W^{(\ell)} N^{(\ell)} + W_d^{(\ell)} \sum_{\text{daughters}} N_d^{(\ell)} \right)$$

Intuition: **Create representation of node considering its neighbours**

- Custom weights for parent (W_p), node (W), daughters (W_d)
- Independent of daughter ordering
- Normalise adjacency matrix
 - Prevent over-representation in high multiplicity decays

Normalised
Laplacian

$$\tilde{A} = A + I_N$$

$$\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$$

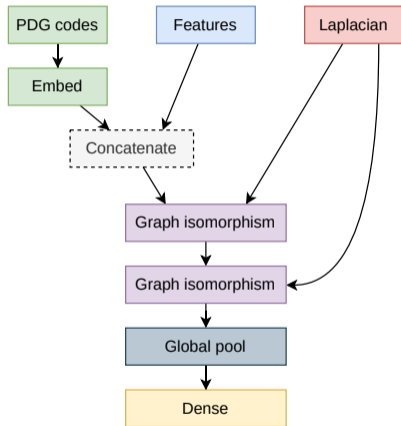
$$\tilde{L} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$$

Special case of:

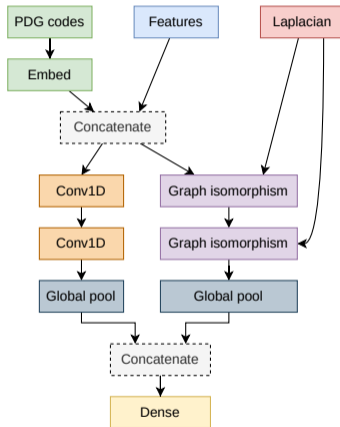
K. Xu, W. Hu, J. Leskovec, S. Jegelka, [How Powerful are Graph Neural Networks?](#) (CoRR 2018)

Training

- Train on 250k events (validate on 10%)
- Test on 50k independent events
- Batch normalisation, dropout, class weights, early stopping, reduce LR on plateau, model checkpoint (save only best), ...

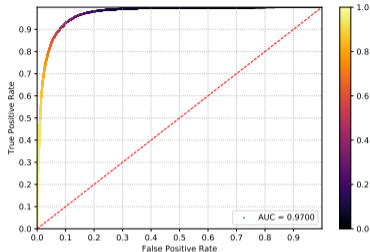


- Train on 250k events (validate on 10%)
- Test on 50k independent events
- Batch normalisation, dropout, class weights, early stopping, reduce LR on plateau, model checkpoint (save only best), ...
- Additional convolutional 1D for full reconstruction dataset

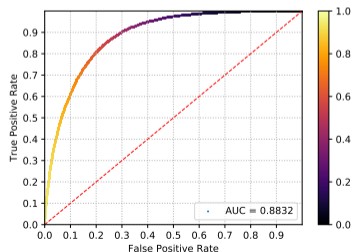


Training

- Train on 250k events (validate on 10%)
- Test on 50k independent events
- Batch normalisation, dropout, class weights, early stopping, reduce LR on plateau, model checkpoint (save only best), ...
- Additional convolutional 1D for full reconstruction dataset



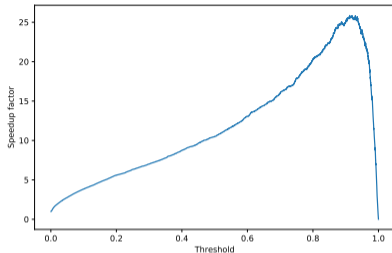
(a) TDCPV



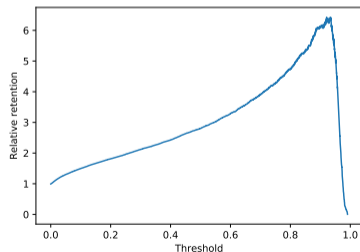
(b) Full reconstruction

Training

- Train on 250k events (validate on 10%)
- Test on 50k independent events
- Batch normalisation, dropout, class weights, early stopping, reduce LR on plateau, model checkpoint (save only best), ...
- Additional convolutional 1D for full reconstruction dataset
- Insert NumPy-based module into Belle II analysis framework for inference



(a) TDCPV

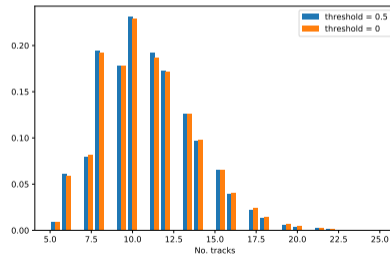
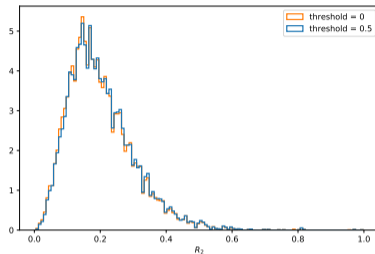
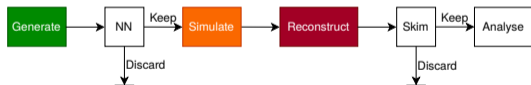


(b) Full reconstruction

Bias check

Compare event-level kinematics:

- Pass skim = True
- Pass skim and NN = True positive



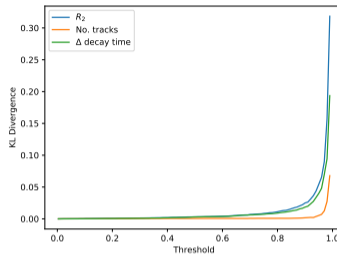
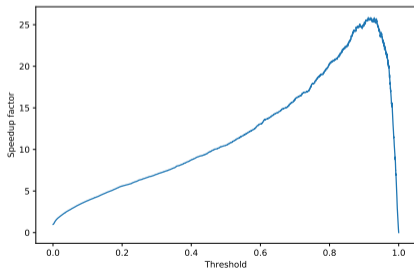
Bias check

Compare event-level kinematics:

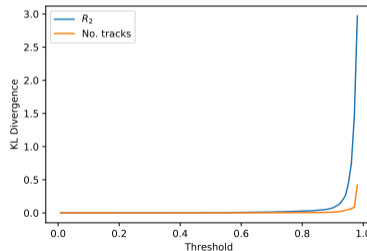
- Pass skim = True
- Pass skim and NN = True positive

Kullback-Leibler divergence of Q from P :

$$D_{\text{KL}}(P \parallel Q) = - \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{Q(x)}{P(x)} \right)$$



(a) TDCPV



(b) Full reconstruction

- Belle II has begun data taking
 - simulation will need to keep up
- Simulations for the full 50 ab^{-1} too computationally expensive
 - Requires smarter solutions
- Propose to use NN to go from: **simulate everything** → **simulate necessary**
 - Must be general enough to handle each physics working-group case
- Shown potential for orders of magnitude speedup and quantification of bias

Current work:

- Scale up datasets and bias checks
- Implement bias mitigation

Thank you

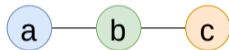
Backup

Original Graph Convolutional Networks (GCN)

Propagation rule of layer activations $H^{(l)}$

$$H^{(l+1)} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right)$$

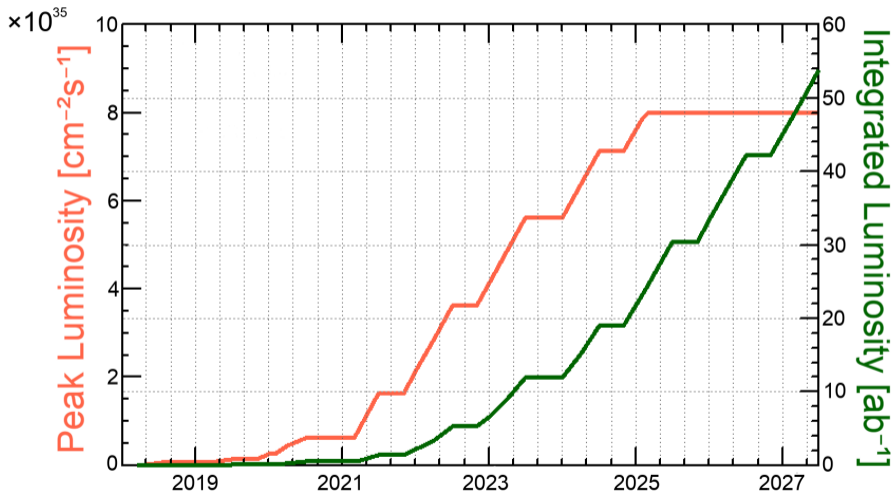
$$\begin{aligned} H^{(0)} &= X \\ \tilde{A} &= A + I_N \\ \tilde{D}_{ii} &= \sum_j \tilde{A}_{ij} \end{aligned}$$



$$\tilde{A}^{N \times N} = A + I = \begin{matrix} & \begin{matrix} a & b & c \end{matrix} \\ \begin{matrix} a \\ b \\ c \end{matrix} & \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix} \end{matrix}$$

Thomas N. Kipf, Max Welling, [Semi-Supervised Classification with Graph Convolutional Networks](#) (ICLR 2017)

Luminosity projection



TDCPV divergence (overload)

