

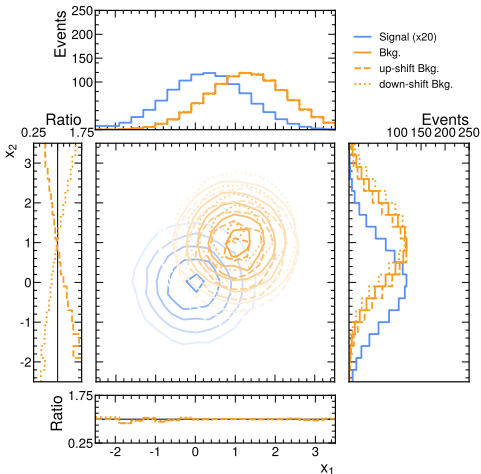
Development of systematic-aware neural network trainings for binned-likelihood-analyses at the LHC

Based on [CMS-PAS-MLG-23-005](#)

11. 09. 2024

Markus Klute, **Artur Monsch**, Lars Sowa, Roger Wolf

Starting point: a synthetic example



Objective: signal strength $r_s \pm \Delta r_s$ in a given model

Common problems for classification tasks (i.e. left):

- High imbalance between process yields
- Processes overlap in some phase space regions
- Processes are (usually) affected by uncertainties

→ **Utilizing NN models for process classification**

Conventional analyses

○○

Δ information

○○○○○○

Effects from added Δ

○○○

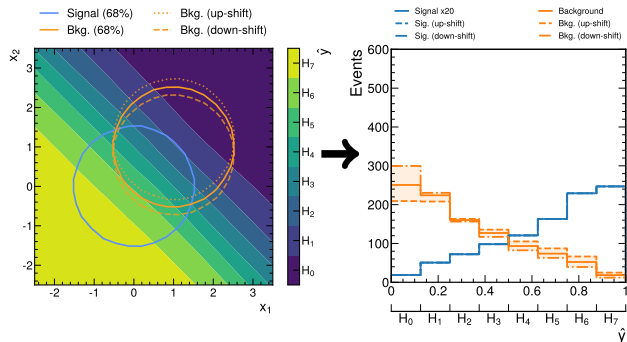
$H(125) \rightarrow \tau\tau$ application

○○○○○○○

Summary

○○○

NN models for conventional process separation



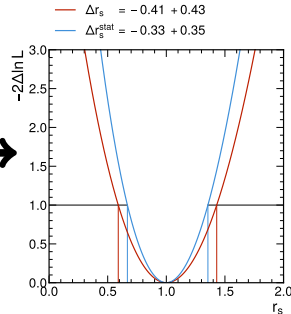
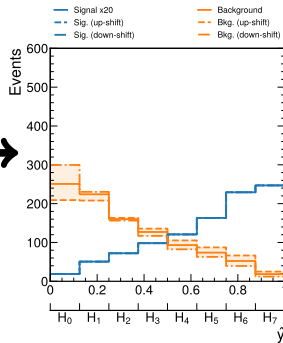
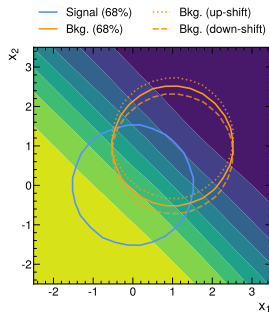
Cross entropy NN training (CENNT)

- Training objective: process separation
- Utilization of nominal signal (S) and background (B) dataset

Resulting NN output

- Nominal S and B datasets
- Systematic variations ΔS_j , ΔB_j for each uncertainty source j

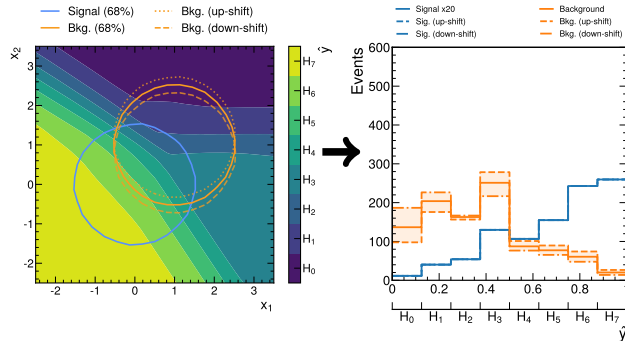
$$r_s \pm \Delta r_s$$



Final results obtained from statistical inference based on binned likelihood

Additional information about systematic variations

NN model for Δr_s minimization



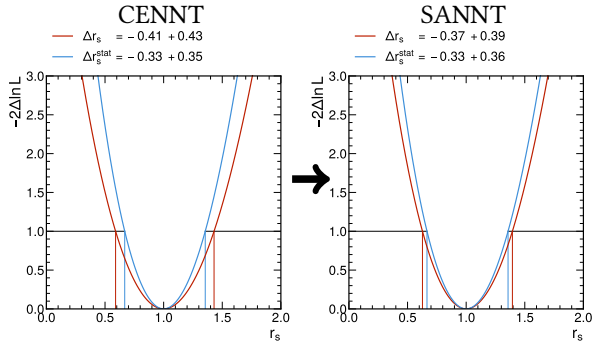
Systematic aware NN training (SANNT)

- Training objective: Δr_s minimization, now aligns with analysis objective
- Not a process separation anymore: High S/B bins creation with low syst. variations

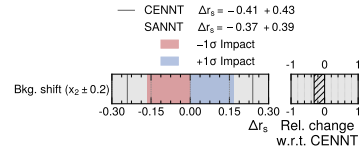
Resulting NN output

- Nominal S and B datasets
- Systematic variations ΔS_j , ΔB_j for each uncertainty source j

Consequences for Δr_s



- Statistical uncertainty comparable
- Redistribution of events reduces impact of uncertainty on r_s
- Reduction of overall uncertainty



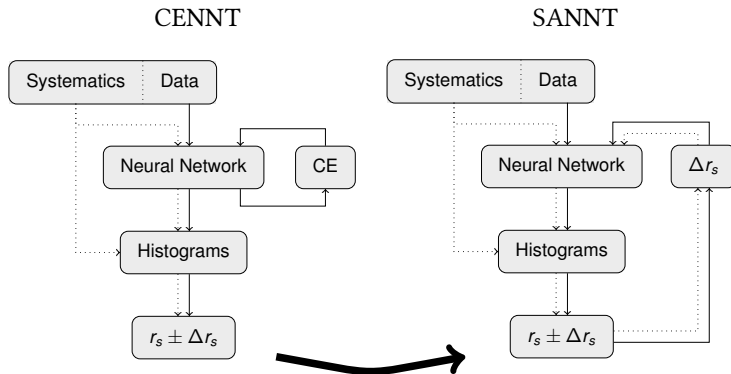
Changes to the training procedure summarized

Use CE pretraining achieving:

- Process separation
- Good starting point for SANNT

Main difference: SANNT vs. CENNT :

- Changed training objective
- Added information about systematic variations to the training



Δr_s as training objective

1 Modification to existing likelihood $\mathcal{L}(\{k_i\}, \{r_s\}, \{\theta_j\})$ ([backup](#)):

- Extend λ_i in $\mathcal{P}(k_i, | \lambda_i,)$ to $\lambda'_i = \sum_s r_s S_{si} + \sum_b B_{bi} + \tilde{\Delta}_i$ with systematic shifts $\tilde{\Delta}_i$, set $\theta_j = 0 \forall j$ and use Asimov dataset, replacing $k_i \rightarrow k_i^A$

$$\tilde{\Delta}_i = \sum_{p \in \{s, b\}} \sum_j \max(0, \theta_{jp}) \left(\Delta_{jp}^{\text{up}} \right)_i + \min(\theta_{jp}, 0) \left(\Delta_{jp}^{\text{down}} \right)_i,$$

- Build computational graph of effects of syst. variations effects for each process $p \in \{s, b\}$
- Including asymmetries at $\theta_k^\pm \rightarrow 0$ for each deviation from nominal

Δr_s as training objective

- 1 Modification to existing likelihood $\mathcal{L}(\{k_i\}, \{r_s\}, \{\theta_j\})$ ([backup](#)):
 - Extend λ_i in $\mathcal{P}(k_i, | \lambda_i,)$ to $\lambda_i' = \sum_s r_s S_{si} + \sum_b B_{bi} + \tilde{\Delta}_i$ with systematic shifts $\tilde{\Delta}_i$, set $\theta_j = 0 \forall j$ and use Asimov dataset, replacing $k_i \rightarrow k_i^A$
- 2 Asymptotically estimate Δr_s from $\sqrt{F_{r_s, r_s}^{-1}}$ using full, weighted training/validation dataset

$$F_{x_i x_j} = \mathbb{E} \left[\frac{\partial^2}{\partial x_i \partial x_j} \left(-\log \mathcal{L} \right) \right]_{x_i x_j \in \{\{r_s\}, \{\theta_j\}\}} = \left(\text{Hess} \left(-\log \mathcal{L} \right) \right)_{x_i x_j \in \{\{r_s\}, \{\theta_j\}\}}$$

For n signals: sum over n diagonal elements of $\sqrt{F_{r_s, r_s}^{-1}}$

Caviat: Backpropagation

Training should have a **non breaking backpropagation**

- Every computation step of Δr_s must be differentiable
- Gradient for histograms is usually not provided (i.e. PYTORCH, TENSORFLOW)

Approach: introduce bin-wise **custom functions** \mathcal{B}_i replacing histogram gradient (adapted from [1])

- Histogram remains unchanged in the forward pass
- Use \mathcal{B}_i only during backward pass
- Choice of \mathcal{B}_i has an effect on training procedure

→ Improvement of \mathcal{B}_i can enable application to more complex tasks

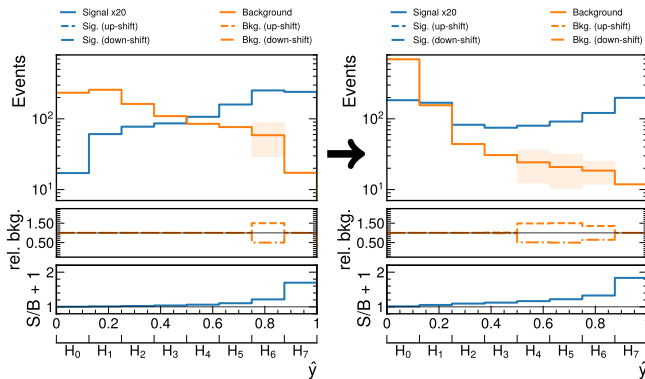
Effects of systematic aware training

Visualizing effects of systematic variations

Injecting uncertainty for B in bin 6 after pretraining (left)

After subsequent SANNT training (right):

- S events
 - Approx. maintained in S -enriched regions
- B events
 - Further removed from S -enriched region mostly to the left in the histogram
- Unct. affected B events moved to the left
 - Effect of unct. is reduced

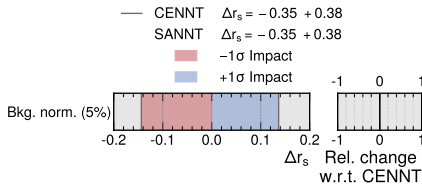
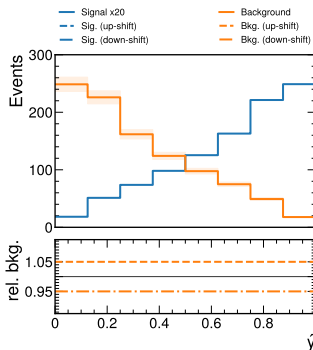


Systematic variations that can not be addressed

- Best process separation achieved after pretraining
 → lowest possible Δr_s^{stat} achieved

In case of pure normalization uncertainties:

- Events are present in all bins and are equally affected
 → no phase space that could lead to improvement of Δr_s w.r.t. Δr_s^{stat}



Systematic variations that should not be considered

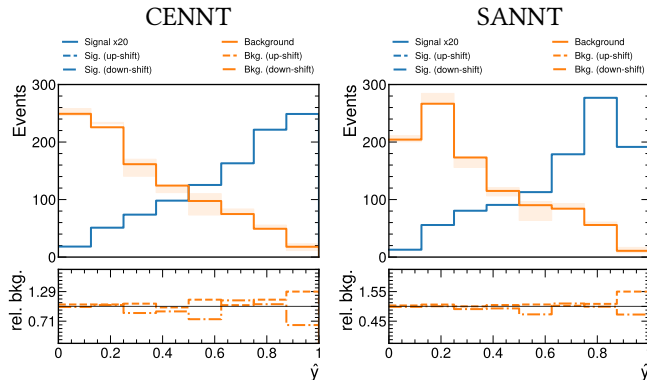
Low statistics uncertainty sample

- **Downsampled** B events from 5% normalization shift (previous slide)

Fluctuating results are retrieved in both cases

SANNT will pick up the fluctuations and try to minimize their effect on $\Delta r_s!$

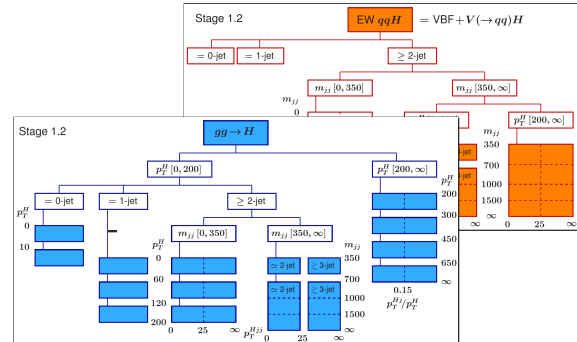
→ **Good uncertainty model description remains essential here**



Application on reduced $H(125) \rightarrow \tau\tau$ analysis

ML-based $H(125) \rightarrow \tau\tau$ analysis (HIG-19-010, [2])

- **Differential cross-section** measurement of $H(125)$ production based on [STXS scheme](#)
- Utilized multiclass classification with **5 background** and up to **15 STXS signal classes**

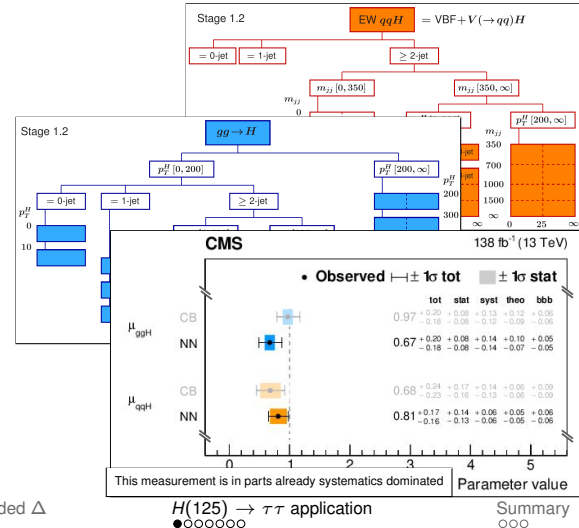


ML-based $H(125) \rightarrow \tau\tau$ analysis (HIG-19-010, [2])

- **Differential cross-section** measurement of $H(125)$ production based on [STXS scheme](#)
- Utilized multiclass classification with **5 background** and up to **15 STXS signal classes**

Future prospects for this and other analyses:

- Statistical uncertainties will decrease (Run3 and HL-LHC)
- Importance of systematic uncertainties will increase
→ Addressing them will become more important



Conventional analyses
○○

Δ information
○○○○○

Effects from added Δ
○○○

Setup overview of application on reduced HIG-19-010 analysis ([3])

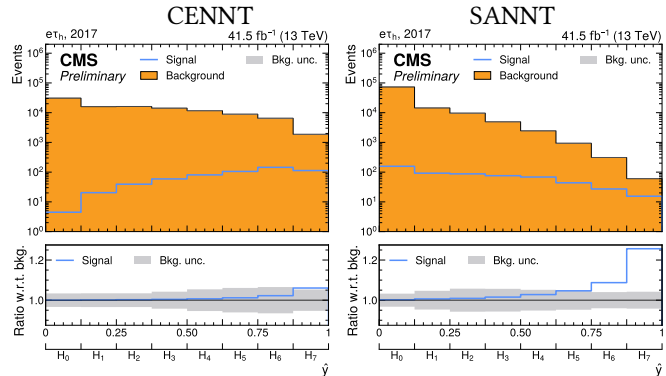
- Using a subset of the dataset used by [2]
 - Final state: $e\mu, \mu\tau_h, e\tau_h, \tau_h\tau_h$
 - Era: 2016, 2017, 2018
- Selecting 86 theoretical and experimental uncertainties
- NN input: 15 variables used as in [2]
- Full-batch training with evaluation based on two-fold cross validation scheme
- Binary classification: all S_s (B_b) processes grouped
- Multiclass classification: 5 B and 2 S processes

Application on reduced analysis of HIG-19-010

Binary classification

CENNT vs. SANNT: Binary classification

- Conceptual differences of SANNT:
 - S/B separation not the primary target
 - No relational bin information for SANNT: bin-wise ordering due to CE pretraining
- Improvements in process separation due to B reduction in S-enriched bins
- Largest total B unct. contributions move away from S-enriched bins



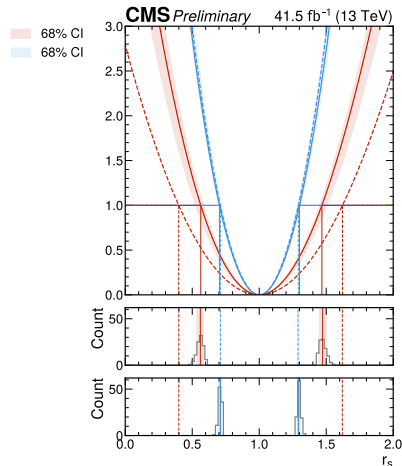
Results of an ensemble test

- Ensemble configuration:
 - Sample size: 100 repetitions
 - Changing NN weight initialization
 - Confidence interval derived from Δr_s

- Fit to Asimov data ($S + B$ model with SM-like signal $r_s = 1$)

- Median expectation of Δr_s w.r.t. CENNT reduction mainly by reducing systematic component of Δr_s

SANNT
 — $\Delta r_s = -0.44 + 0.47$
 — $\Delta r_s^{\text{stat}} = -0.30 + 0.30$
 CENNT
 - - - $\Delta r_s = -0.60 + 0.62$
 - - - $\Delta r_s^{\text{stat}} = -0.29 + 0.29$



Conventional analyses
 ○○

Δ information
 ○○○○○○

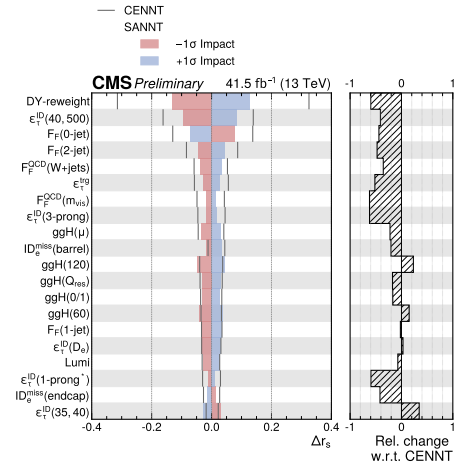
Effects from added Δ
 ○○○

$H(125) \rightarrow \tau\tau$ application
 ○○○●○○○

Summary
 ○○○

Comparison of uncertainty impacts

- Importance ordering as obtained after CENNT
(Detailed description: [backup](#))
- Reduction of uncertainties with largest impact on Δr_s
→ Comparable reduction of Δr_s can be achieved using only a subset of uncertainties with largest impact on Δr_s
- Note: Normalization uncertainties in binary classification can show a shape-changing effect if applied on a subset process of S, B



Application on reduced analysis of HIG-19-010

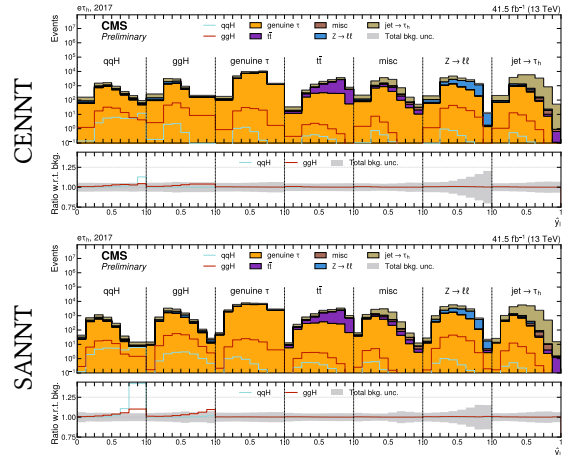
Multiclass classification

SANNT: Multiclass classification

- Considering uncertainty for two signal processes
 - Gluon fusion (ggH)
 - Vector boson fusion (qqH)
- No relational information to predefined classes nor bins
- Sustain concept of output classes

$$L_{\text{SANNT}}^{\text{mult.}} = \sum_s \Delta r_s + \omega_\lambda g(\cdot), \text{ with}$$

- $g(\cdot) = \max(\text{CE}' - \text{CE}'_{\min}{}^{\text{pretrain.}}, 0)$
- and learnable ω_λ [4]



Uncertainty reduction for SANNT multiclass classification

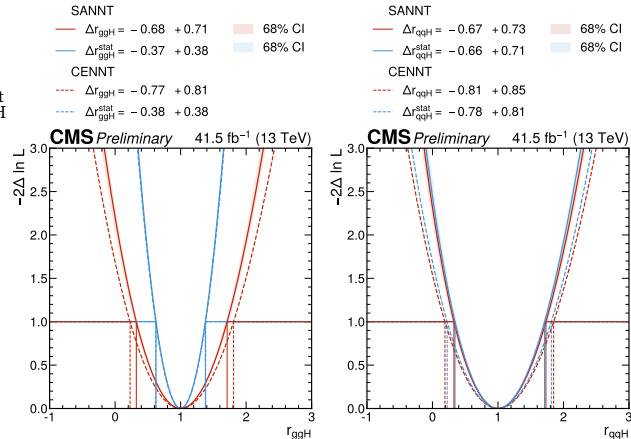
qqH has low process yield in comparison to ggH

- Δr_{qqH} still statistically dominated
- Reduction of Δr_{qqH} through minimization of Δr_{qqH}^{stat}
→ Improvement in process separation

ggH is more distributed across multiple classes

- Process enrichment in few bins more difficult
→ only minor improvement in Δr_{ggH}^{stat}
- Major reduction of Δr_{ggH} through reduction of systematic uncertainty contribution

r_{inc} : [backup](#)



Conventional analyses
○○

Δ information
○○○○○

Effects from added Δ
○○○

$H(125) \rightarrow \tau\tau$ application
○○○○○●

Summary
○○○

Summary

Summary

- Fundamental SANNT studies:
 - General idea of functionality through synthetic examples
 - Effects of various uncertainty sources on training outcomes
 - Effects of poorly understood uncertainty models, regardless of applied method

- Application of SANNT:
 - Systematic NN training is applicable to real-life analysis tasks with considerably higher complexity
 - Extension to multiple classes while maintaining interpretability is possible
 - Significant sensitivity improvement over classical CE-based training using 86 uncertainties

References I

- [1] Stefan Wunsch et al. “Optimal Statistical Inference in the Presence of Systematic Uncertainties Using Neural Network Optimization Based on Binned Poisson Likelihoods with Nuisance Parameters”. In: *Computing and Software for Big Science* 5.1 (Jan. 2021), p. 4. ISSN: 2510-2044. DOI: [10.1007/s41781-020-00049-5](https://doi.org/10.1007/s41781-020-00049-5). URL: <https://doi.org/10.1007/s41781-020-00049-5>.
- [2] Armen Tumasyan et al. “Measurements of Higgs boson production in the decay channel with a pair of τ leptons in proton–proton collisions at $\sqrt{s} = 13$ TeV”. In: *Eur. Phys. J. C* 83.7 (2023). All the figures and tables, including additional supplementary figures and tables, can be found at <http://cms-results.web.cern.ch/cms-results/public-results/publications/HIG-19-010> (CMS Public Pages), p. 562. DOI: [10.1140/epjc/s10052-023-11452-8](https://doi.org/10.1140/epjc/s10052-023-11452-8). arXiv: [2204.12957](https://arxiv.org/abs/2204.12957). URL: <https://cds.cern.ch/record/2807752>.
- [3] Development of systematic-aware neural network trainings for binned-likelihood-analyses at the LHC. Tech. rep. Geneva: CERN, 2024. URL: <http://cds.cern.ch/record/2905411>.

References II

- [4] John Platt and Alan Barr. “Constrained Differential Optimization”. In: Neural Information Processing Systems. Ed. by D. Anderson. Vol. 0. American Institute of Physics, 1987.
URL:
<https://proceedings.neurips.cc/paper/1987/file/a87ff679a2f3e71d9181a67b7542122c-Paper.pdf>.

Backup

Likelihood formulation

- All histogram bins (i) of all classes (c) enter as inputs to an extended binned likelihood function \mathcal{L} , including all nuisance parameters (θ_j)

$$\mathcal{L}(\{k_i\}, \{r_s\}, \{\theta_j\}) = \prod_c \left[\prod_i \mathcal{P}(k_i | \lambda_i) \prod_j \mathcal{C}_j(\tilde{\theta}_j | \theta_j) \right]_c$$

$$\lambda_i = \sum_s r_s \mathcal{S}_{si}(\{\theta_j\}) + \sum_b B_{bi}(\{\theta_j\})$$

with

- Poisson distribution $\mathcal{P}(\cdot | \cdot)$ constructed from k_i observed and λ_i expected events per bin i
 - Use $S + B$ model with scaling parameter(s) r_s
 - Uncertainties are taken into account in the form of nuisance parameters θ_j , following predefined pdf's $\mathcal{C}_j(\cdot | \cdot)$
- Estimation of the uncertainty Δr_s on r_s through Asimov dataset $D_H^A: k_i \rightarrow k_i^A$

Choice of \mathcal{B}_i

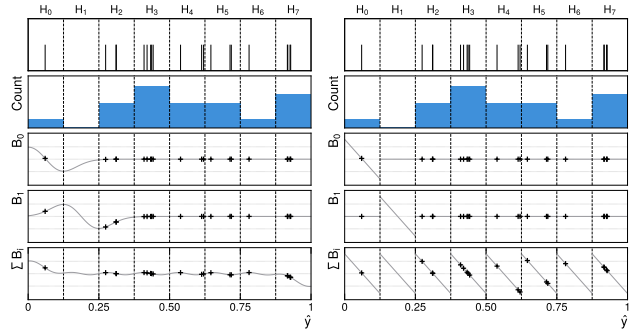
From [1] (left): Derivative of Gaussian PDF

- Low Gradient amplitude at center of bins
- Observing undesired concentration of events in very few bins in $H(\hat{y})$

→ Hard to apply to more complex tasks

Modification (right): linear function within bin boundaries

- Maintains injectivity
- Restricts range to corresponding bin H_i
- Keeps bins indistinguishable



Event evolution during training

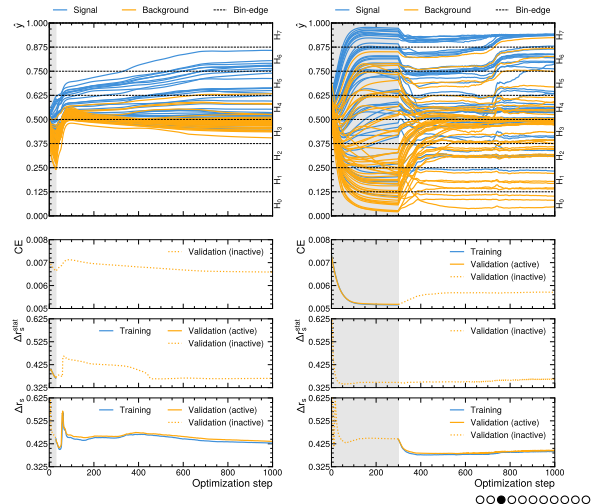
Using \mathcal{B}_i from [1] with same setup ($\Delta x_2 = \pm 1$)

- Strong concentration in \hat{y}
- Independent from [choice of pretraining](#)

Pretraining change to CE:

- Equivalent to Δr_s^{stat} minimization
- Providing better starting point for minimization of systematic component of Δr_s

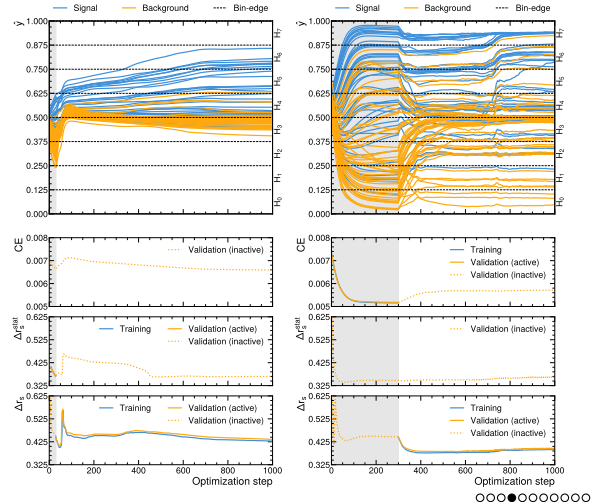
Identity operation (STE) for \mathcal{B}_i : [similar](#) behaviour to [1]



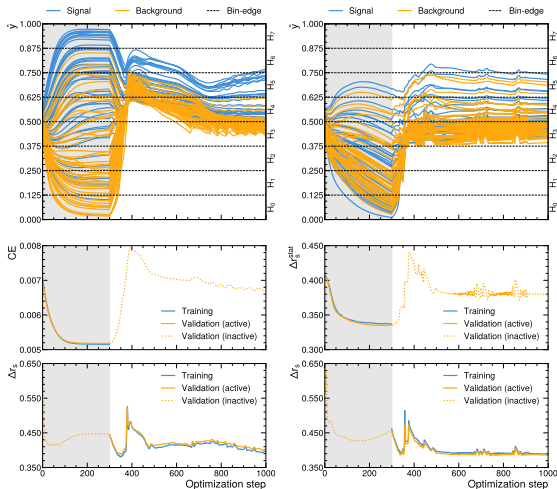
Choice of \mathcal{B}_i : evolution during training

Difference between CE and $\Delta r_S^{(stat)}$

Due to the product in $\mathcal{L} = \prod_i \mathcal{P}(k_i | \lambda_i)$ there is no sequential bin information when minimizing $\Delta r_S^{(stat)}$



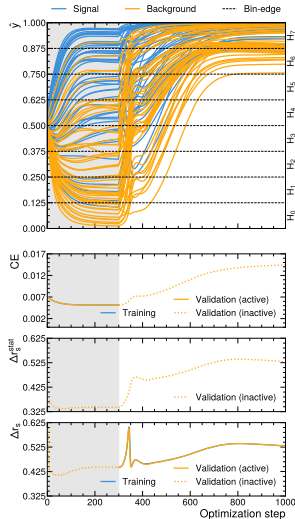
Choice of \mathcal{B}_i : Evolution during training - different pretraining



Derivative of Gaussian PDF as choice for \mathcal{B}_i

- Collapse of \hat{y} into single bins independent of pretraining loss (CE or Δr_s^{stat}) or duration (here 300 optimization steps)
- Collapse still prominent, \hat{y} is less spread after pretraining/with shorter pretraining duration
- Best result achieved here (left) for 300 optimization steps of CE as pretraining: during collapse phase

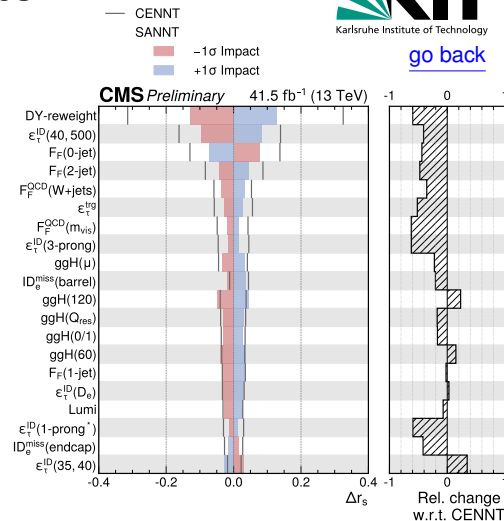
Evolution using Identity operation (STE) as \mathcal{B}_i



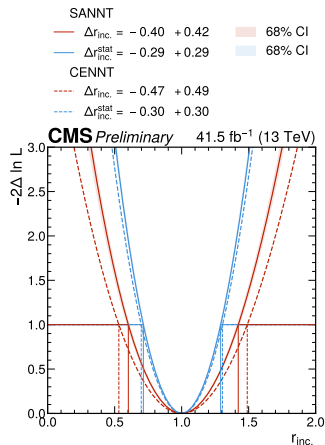
- Using identity for \mathcal{B}_i leads to same problem as \mathcal{B}_i as proposed in [1]
- Bins where the collapse into single bins occurs vary, depending on the seed used for weight initialization

Description of most impactful uncertainties (binary classification)

Label	Type	Process	Rank	Norm	Shape	Comment
$\epsilon_{\tau}^{\text{trig}}$	τ -Trigger	EMB	6	-	✓	-
$\epsilon_{\tau}^{\text{ID}}(D_e)$	τ -ID	MC, EMB	16	✓	-	Discr. against e
$\epsilon_{\tau}^{\text{ID}}(35, 40)$	τ -ID	EMB	20	-	✓	$35 < p_T^{\text{th}} < 40$ GeV
$\epsilon_{\tau}^{\text{ID}}(40, 500)$	τ -ID	EMB	2	-	✓	$40 < p_T^{\text{th}} < 500$ GeV
$\epsilon_{\tau}^{\text{ID}}(1\text{-prong}^*)$	τ -ID	EMB	18	-	✓	One $\pi^+ + \pi^0$'s
$\epsilon_{\tau}^{\text{ID}}(3\text{-prong})$	τ -ID	EMB	8	-	✓	Three π^+ 's
$F_F(0\text{-jet})$	Norm.	F_F	3	-	✓	$N_{\text{jet}} = 0$
$F_F(1\text{-jet})$	Norm.	F_F	15	-	✓	$N_{\text{jet}} = 1$
$F_F(2\text{-jet})$	Norm.	F_F	4	-	✓	$N_{\text{jet}} = 2$
$F_F^{\text{QCD}}(m_{\text{vis}})$	Non-closure	F_F	7	-	✓	In m_{vis}
$F_F^{\text{QCD}}(W+\text{jets})$	Subtr.	F_F	5	-	✓	Subtr. of MC
ggH(μ)	Theory	ggH	9	-	✓	μ_r and μ_f
ggH(Q_{res})	Theory	ggH	12	-	✓	Resummation
ggH(0/1)	Theory	ggH	13	-	✓	0 \rightarrow 1 jet migr.
ggH(60)	Theory	ggH	14	-	✓	p_T^{H} migr.
ggH(120)	Theory	ggH	11	-	✓	p_T^{H} migr.
ID $_e^{\text{miss}}$ (barrel)	e-miss-ID	MC	10	-	✓	Barrel
ID $_e^{\text{miss}}$ (endcap)	e-miss-ID	MC	19	-	✓	Endcap
DY-reweight	Reweight	MC	1	-	✓	In $p_T^{\mu\mu}$ and $m_{\mu\mu}$
Lumi	Luminosity	MC	17	✓	-	-



Uncertainty reduction for multiclass classification $r_{inc.}$



Differences to CENNT: In multiclass classification case

Effects on NN output

- Conceptually:

- No relational information about predefined classes (nor bins)
- Given Δr_s or $\sum_s \Delta r_s$ no penalty for events that swap predefined classes
→ may lead to emergence of empty classes (less interpretable)

→ Introduce a penalty term modifying the loss to: $L_{\text{SANNT}}^{\text{mult.}} = \sum_s \Delta r_s + \omega_\lambda g(\cdot)$ as introduced in [4]

- Aim is to sustain the concept of output classes
- $g(\cdot) = \text{CE}' - \text{CE}'_{\min}$ ensures a class assignment, that corresponds to CE'_{\min} as obtained from pretraining
- ω_λ updated at each optimization step (using backpropagation) and set to 0 if $g(\cdot) < 0$
(improvement w.r.t. CE'_{\min})

- Technical differences:

- Different activation function (Sigmoid vs. Softmax)
- In case of SANNT: modified (binary) CE' for pretraining and as constraint

Related work

Estimation of statistical quantities from binned variables (in HEP) is often necessary due to

- Computational cost
- Availability of corrections/shifts

Transition to NN application requires differentiable histogram operation.

The two solutions to which this work is most comparable are (also mentioned in the paper):

INFERN0¹: Histogram approximation through tunable softmax operator

- Training objective: Diagonal element of inverse Fisher-Information matrix (Δr_s)
- Uses mini-batches during training and explicit sampling of F_{x_i, x_j}

neos²: Histogram approximation through kernel density estimation

- Training objective: likelihood ratio
- Uses mini-batches during training

¹DOI: [10.1016/j.cpc.2019.06.007](https://doi.org/10.1016/j.cpc.2019.06.007)

²DOI: [10.1088/1742-6596/2438/1/012105](https://doi.org/10.1088/1742-6596/2438/1/012105)

Downsampled normalization shift

