



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

ML4Jets 2023

Quark/gluon tagging in CMS Open Data with CWoLa and TopicFlow

Ayodele Ore

In collaboration with Matthew J. Dolan and John Gargalionis

Weakly-supervised Q/G tagging

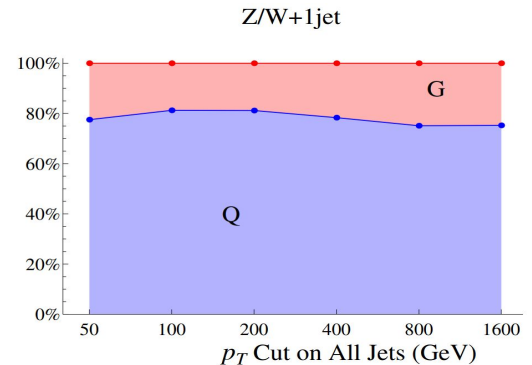
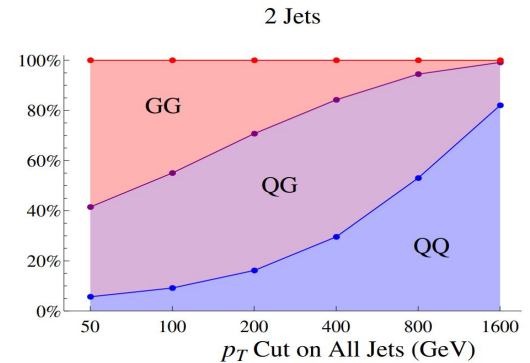
- Fully-supervised learning not for Q/G ideal since:
 - ! Discrimination is sensitive to non-perturbative effects with large uncertainties in MC
 - ! Parton labels not well defined at detector level
- Instead, train on *mixed* samples – obtainable from exp.

$$p_{M_1}(x) = f_1 p_Q(x) + (1 - f_1) p_G(x)$$

$$p_{M_2}(x) = f_2 p_Q(x) + (1 - f_2) p_G(x)$$

- Same optimal classifier for M_1 vs M_2 as Q vs G (CWoLa)
- CMS Open Data is a great testing ground

[JHEP10(2017)174]



[JHEP10(2011)103]

Jet Topics: Disentangled distributions

[Phys.Rev.Lett.120,241602]

- If the mixture fractions are known, the pure distributions can be recovered:

$$\begin{array}{l} p_{M_1}(x) = f_1 p_Q(x) + (1 - f_1) p_G(x) \\ p_{M_2}(x) = f_2 p_Q(x) + (1 - f_2) p_G(x) \end{array} \xrightarrow{\text{Invert}} \begin{array}{l} p_Q(x) = \frac{(1 - f_2) p_{M_1}(x) - (1 - f_1) p_{M_2}(x)}{f_1 - f_2} \\ p_G(x) = \frac{f_1 p_{M_2}(x) - f_2 p_{M_1}(x)}{f_1 - f_2} \end{array} \longrightarrow \text{ROC curve from } \varepsilon_Q \varepsilon_G$$

Jet Topics: Disentangled distributions

[Phys.Rev.Lett.120,241602]

- If the mixture fractions are known, the pure distributions can be recovered:

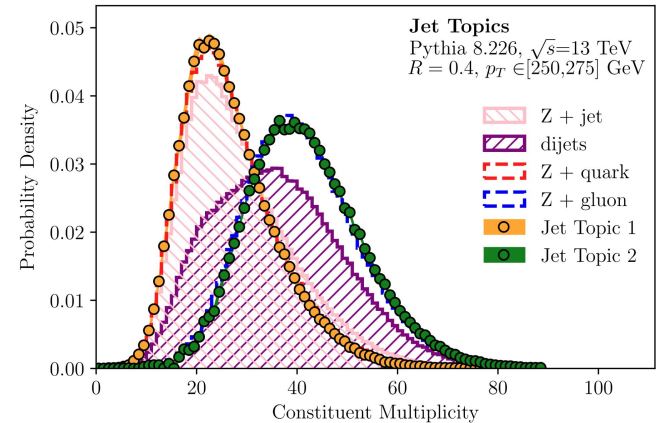
$$\begin{aligned}
 p_{M_1}(x) &= f_1 p_Q(x) + (1 - f_1) p_G(x) \\
 p_{M_2}(x) &= f_2 p_Q(x) + (1 - f_2) p_G(x)
 \end{aligned}
 \xrightarrow{\text{Invert}}
 \begin{aligned}
 p_Q(x) &= \frac{(1 - f_2) p_{M_1}(x) - (1 - f_1) p_{M_2}(x)}{f_1 - f_2} \\
 p_G(x) &= \frac{f_1 p_{M_2}(x) - f_2 p_{M_1}(x)}{f_1 - f_2}
 \end{aligned}
 \xrightarrow{\text{ROC curve from } \varepsilon_Q, \varepsilon_G}$$

- Define “reducibility factors”: $\kappa_{ij} \equiv \min_x \frac{p_{M_i}(x)}{p_{M_j}(x)}$
- If $\kappa_{QG} = \kappa_{GQ} = 0$ (mutual irreducibility) then:

$$f_1 = \frac{1 - \kappa_{12}}{1 - \kappa_{12}\kappa_{21}} \quad f_2 = \kappa_{21}f_1 \quad \leftarrow \text{Measure directly from } M_1, M_2!$$

... Otherwise you need to know κ_{QG} or κ_{GQ}

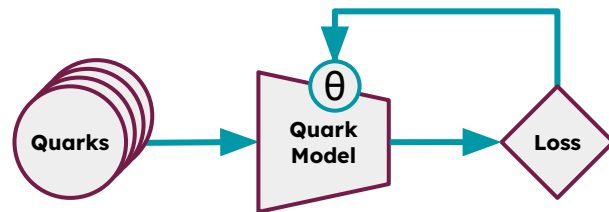
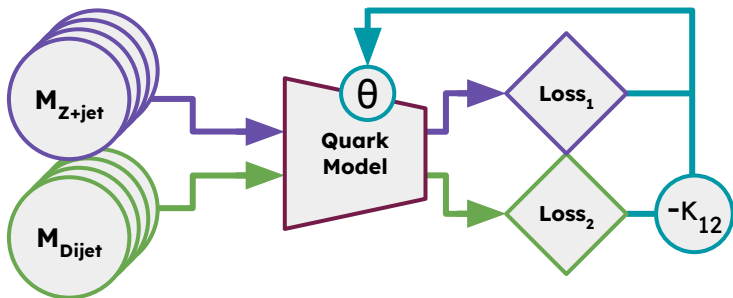
Estimate from Theory/MC



TopicFlow

[AO & Matthew Dolan, [PhysRevD.107.114003](https://arxiv.org/abs/1703.07501)]

- Generative models can be trained to learn topic distributions, given fractions.
 - ✓ Can apply in many dimensions
 - ✓ Can smooth statistics with oversampling
 - ✓ Can access quark/gluon likelihoods (with normalizing flow)



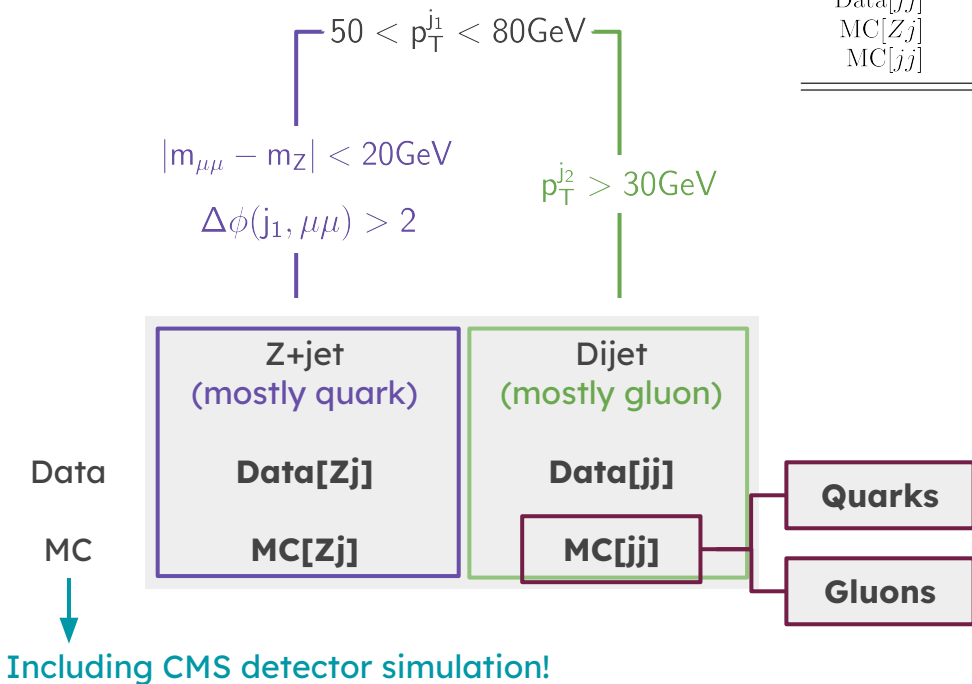
Typically requires pure sample!

$$\text{LOSS}_{\text{Quark}} = \langle L_{\theta}(x) \rangle_{x \sim p_Q(x)}$$

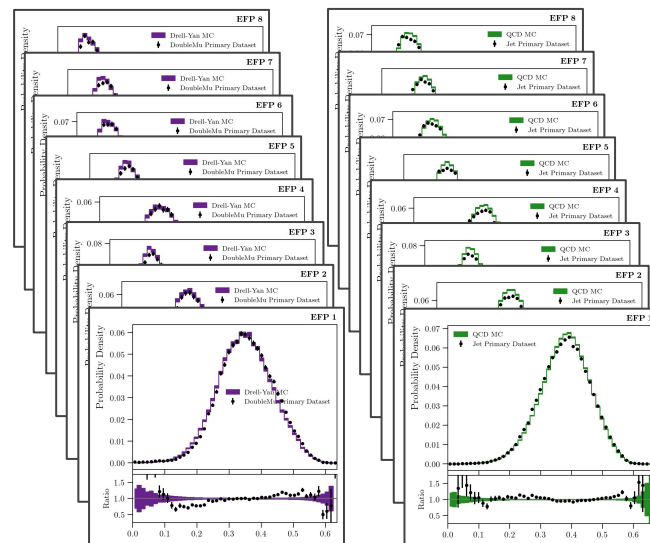
$$= \underbrace{\langle L_{\theta}(x) \rangle_{x \sim p_{M_1}(x)}}_{\downarrow M_1} - \kappa_{12} \underbrace{\langle L_{\theta}(x) \rangle_{x \sim p_{M_2}(x)}}_{\uparrow M_2}$$

CMS Open Data

- 2011 data at @ 7 TeV

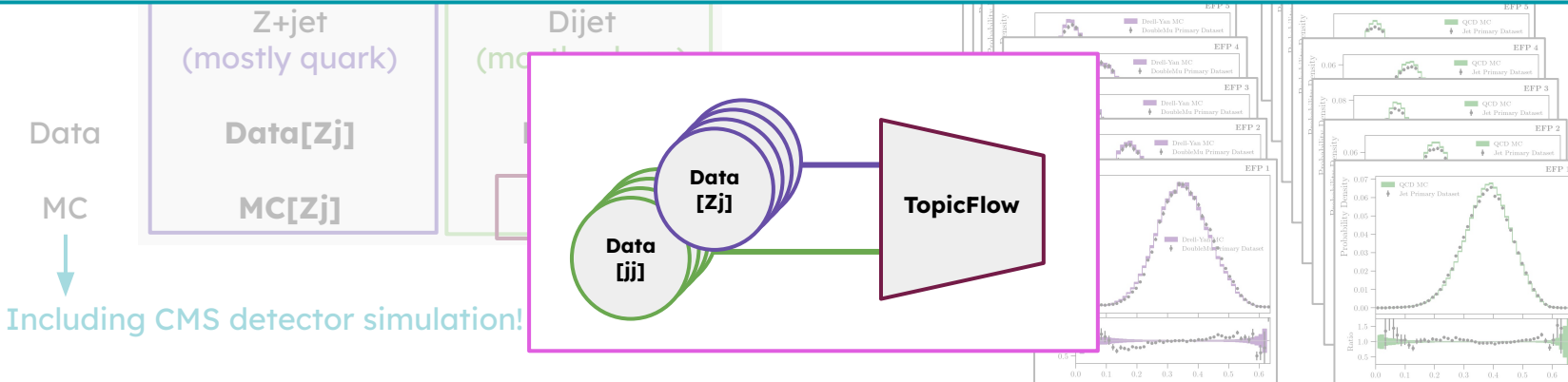
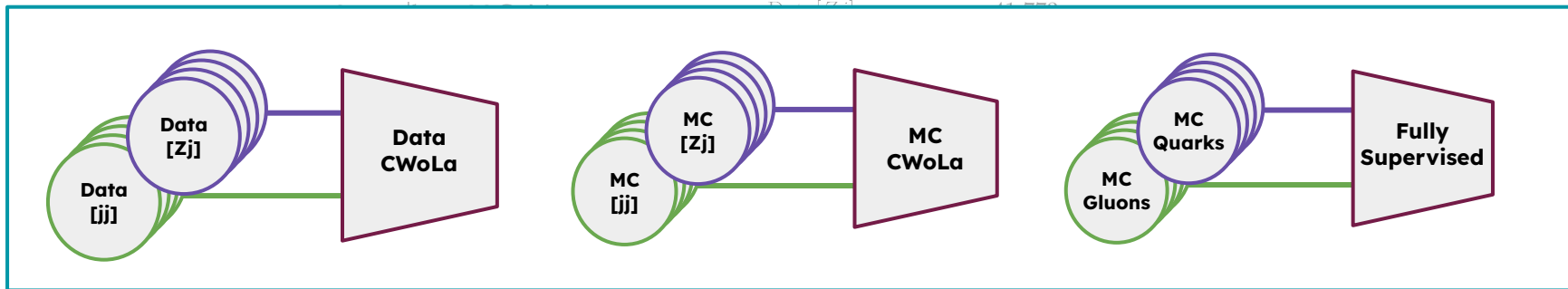


Dataset	Total events	Quarks	Gluons	
Data[Zj]	41,773	-	-	
Data[jj]	82,162	-	-	
MC[Zj]	95,324	70,568	24,756	~74% quarks
MC[jj]	3,064,713	868,556	2,196,157	~30% quarks

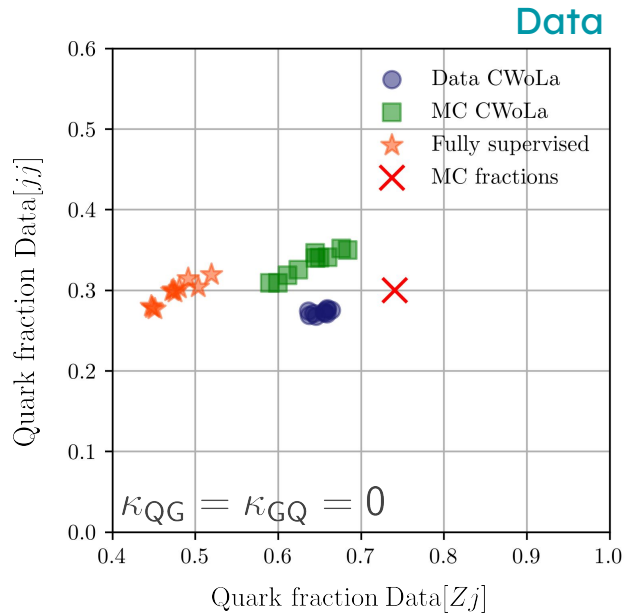


CMS Open Data

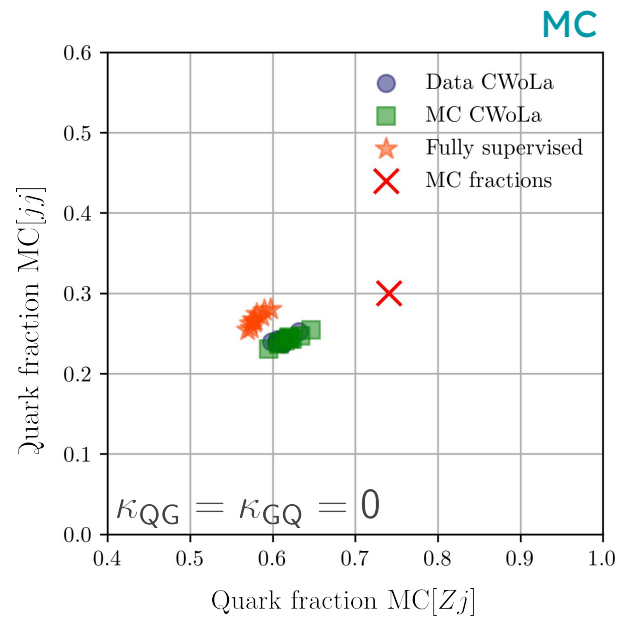
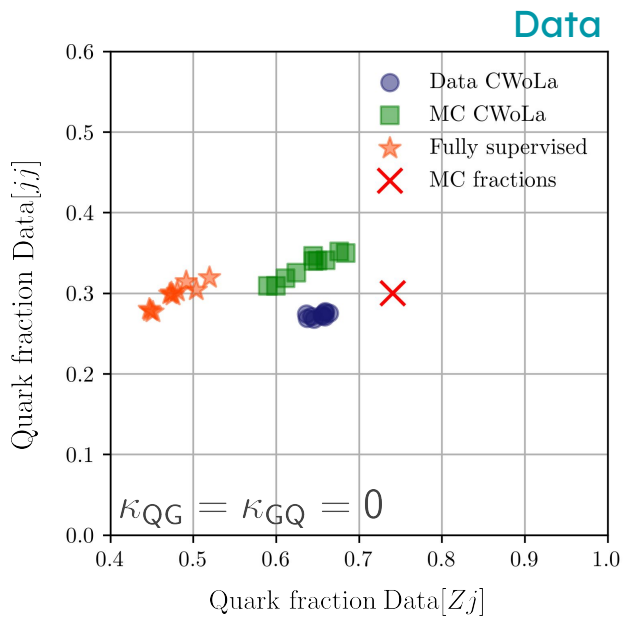
Dataset	Total events	Quarks	Gluons
Di-jet	1.1M	10%	90%
Z+jet	1.1M	90%	10%



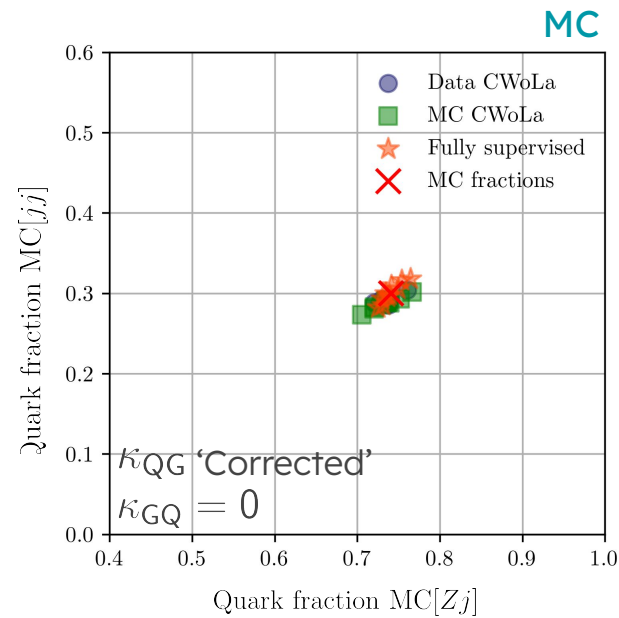
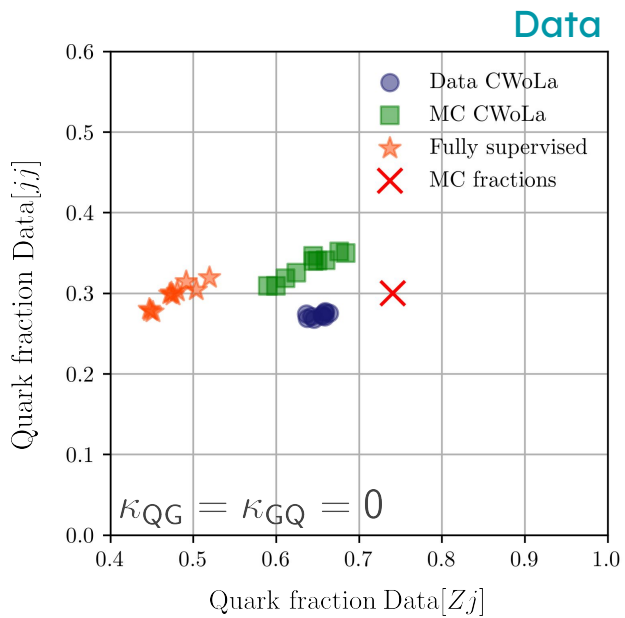
Quark fractions



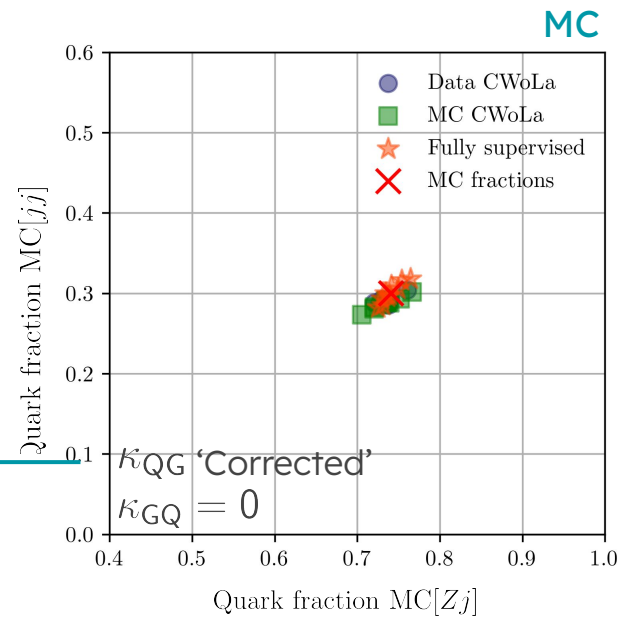
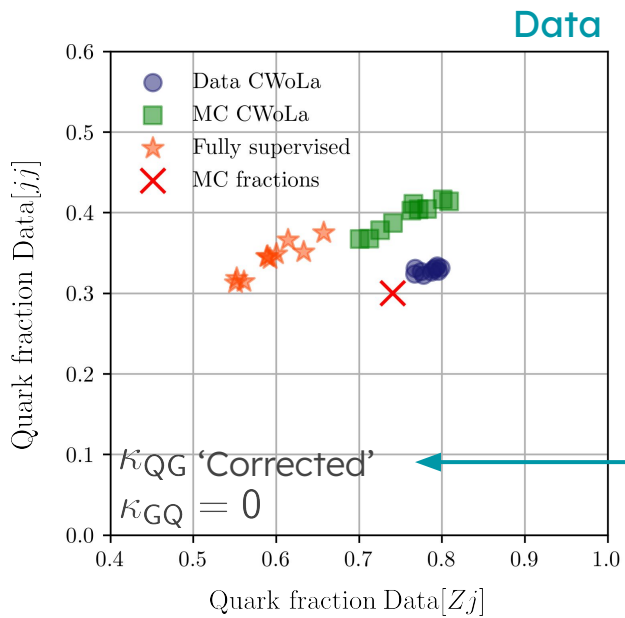
Quark fractions



Quark fractions



Quark fractions

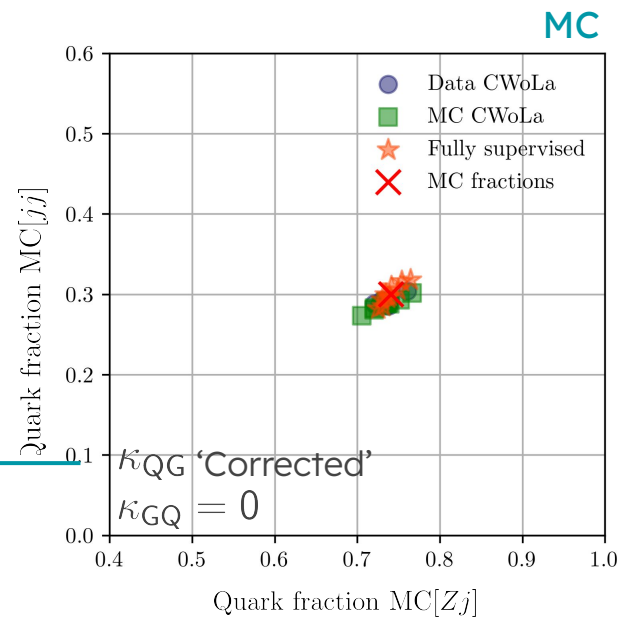
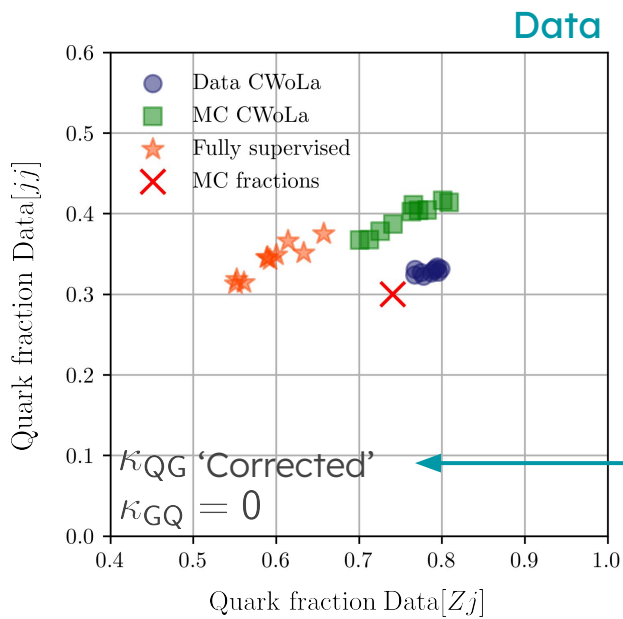


Quark fractions

From Data CWoLa:

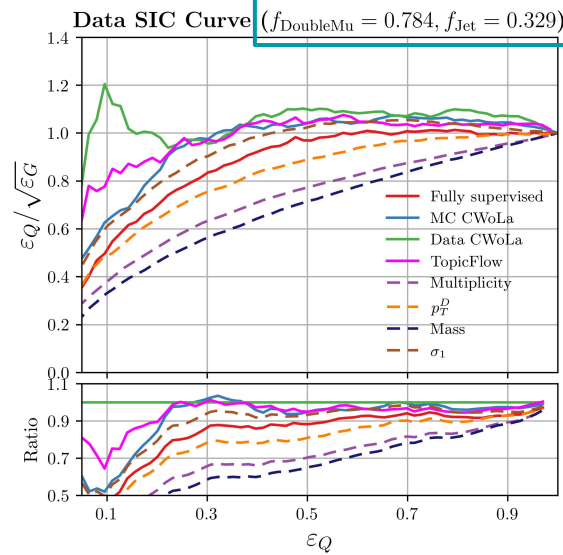
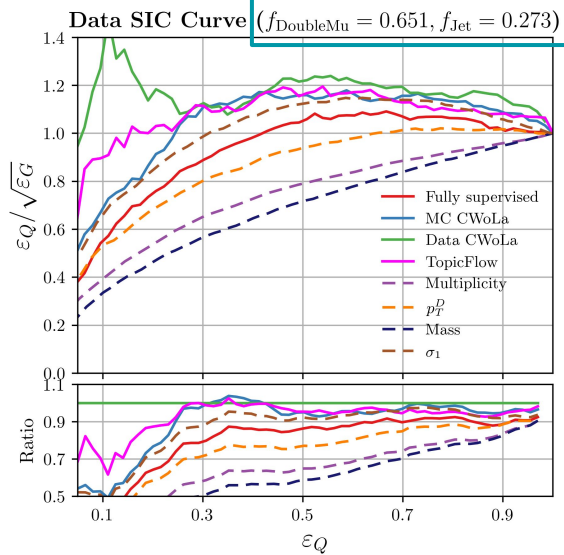
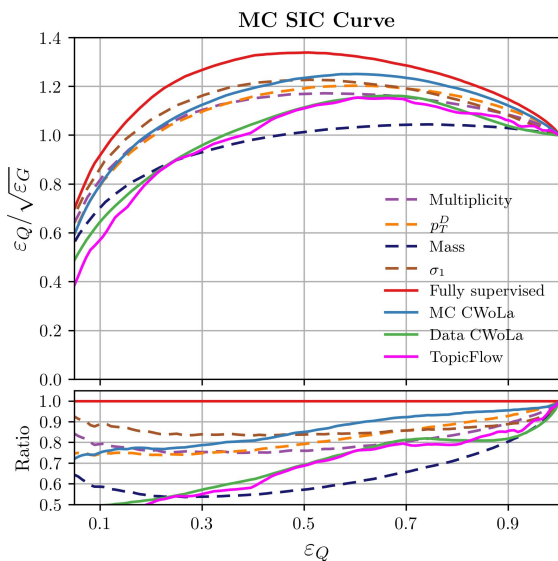
$$\kappa_{QG} = 0 : f_1 = 0.651, f_2 = 0.273$$

$$\kappa_{QG} \text{ 'Corrected' } : f_1 = 0.784, f_2 = 0.329$$



Tagging performance

- Recall estimate for efficiencies on data: $\varepsilon_G(t) = \frac{f_1 \varepsilon_{M_2}(t) - f_2 \varepsilon_{M_1}(t)}{f_1 - f_2}$ $\varepsilon_Q(t) = \frac{(1 - f_2) \varepsilon_{M_1}(t) - (1 - f_1) \varepsilon_{M_2}(t)}{f_1 - f_2}$



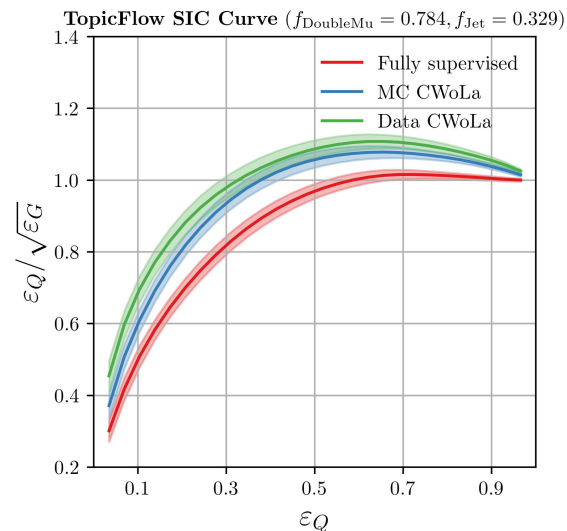
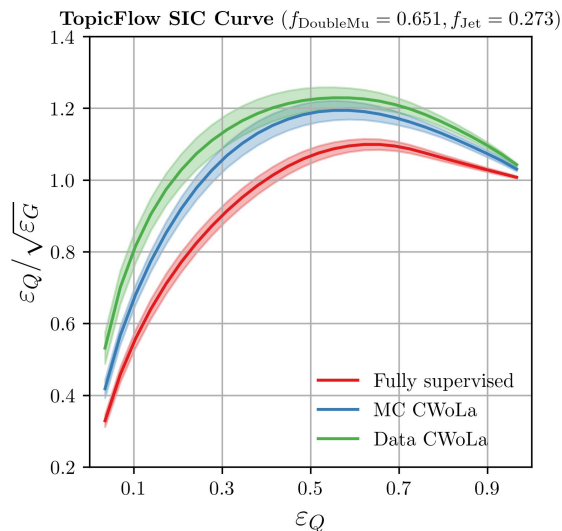
Smoothed ROCs with TopicFlow

$$\varepsilon_Q(t) = \langle \Theta(h(x) - t) \rangle_{x \sim p_Q}$$

$$\varepsilon_G(t) = \langle \Theta(h(x) - t) \rangle_{x \sim p_G}$$

- TopicFlow lets us avoid subtraction and evaluate ROC on samples:

- Error bars from ensemble: (Different distributions compatible with the test dataset)
- Bands smaller than impact of assumed fractions



Summary

- Fully-supervised learning causes train-test domain shift for jet tagging.
- Weakly-supervised learning facilitates training on data:

Classification

CWoLa outperforms MC-supervised models in CMS data.

Classifier rankings unaffected by estimated fraction.

- Future directions:

★ Model/quantify sample dependence

Generative Modelling

TopicFlow can learn pure quark/gluon distributions.

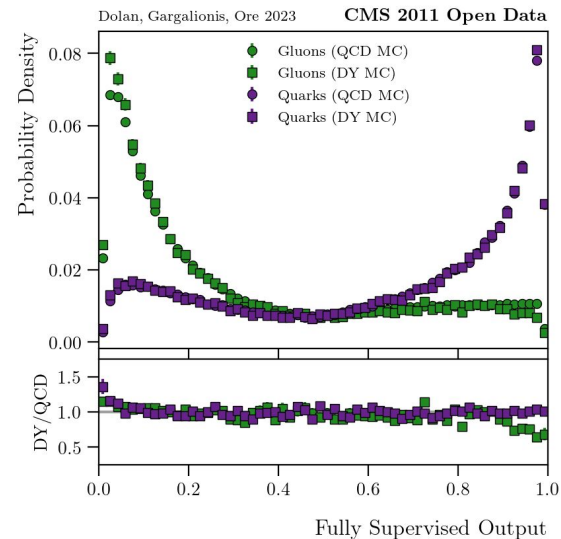
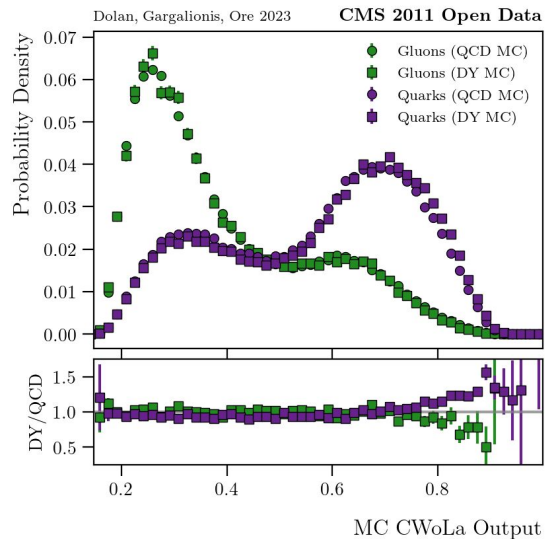
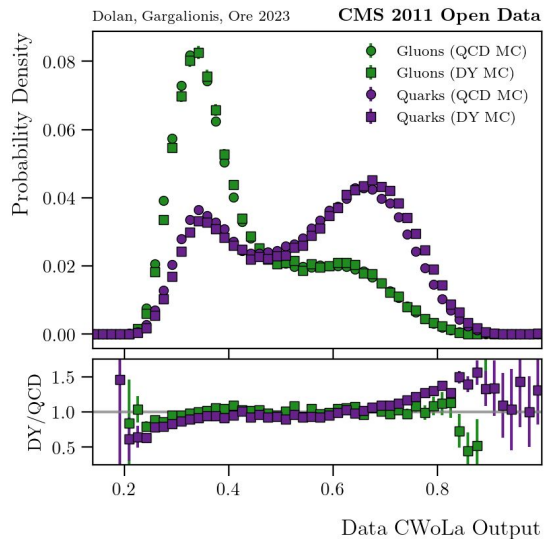
Generative classification competitive with CWoLa.

Oversampling smooths ROC curves.

★ Add f_1, f_2 to TopicFlow optimization

Backup Slides

Backup: Sample independence



Backup: All Jet Topics formulas

$$p_{M_1}(x) = f_1 p_Q(x) + (1 - f_1) p_G(x)$$

$$p_{M_2}(x) = f_2 p_Q(x) + (1 - f_2) p_G(x)$$



$$p_G(x) = \frac{f_1 p_{M_2}(x) - f_2 p_{M_1}(x)}{f_1 - f_2}$$

$$p_Q(x) = \frac{(1 - f_2) p_{M_1}(x) - (1 - f_1) p_{M_2}(x)}{f_1 - f_2}$$

Define “reducibility factors” $\kappa_{ij} \equiv \min_x \frac{p_{M_i}(x)}{p_{M_j}(x)}$

Measure directly from M_1, M_2

If $\kappa_{QG} = \kappa_{GQ} = 0$ (“mutual irreducibility”) then: $f_1 = \frac{1 - \kappa_{12}}{1 - \kappa_{12}\kappa_{21}}$, $f_2 = \kappa_{21}f_1$



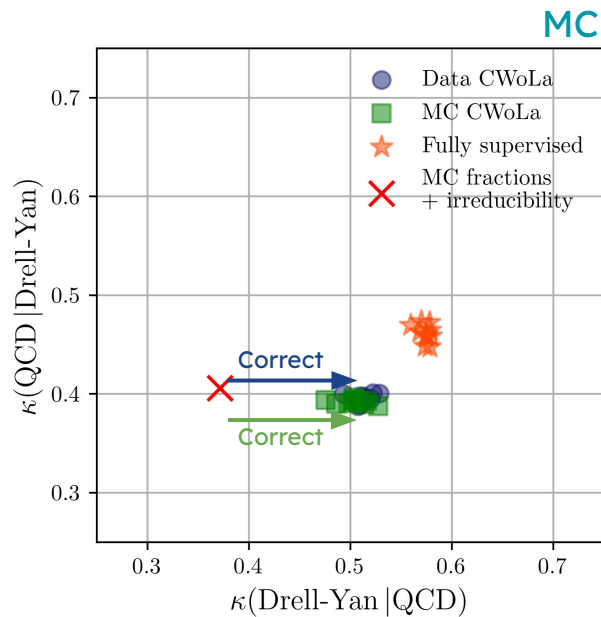
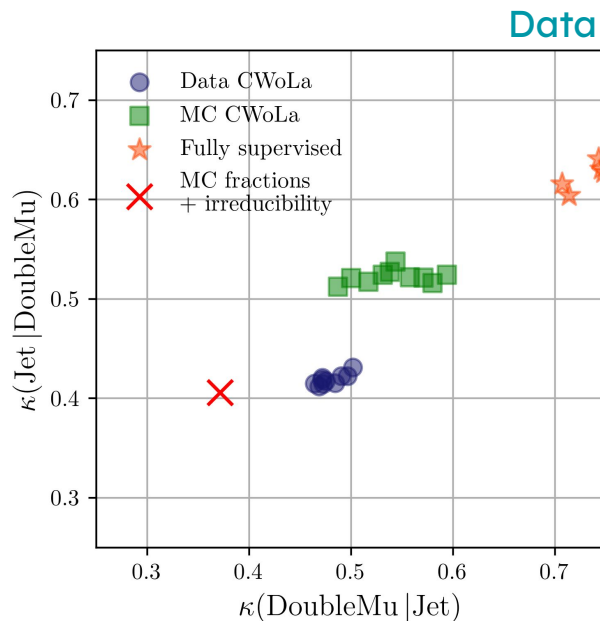
For non-zero κ_{QG} or κ_{GQ} :

← Calibrate with MC/Theory

$$f_1 = \frac{1}{1 - \kappa_{12}\kappa_{21}} \left(\frac{1 - \kappa_{12}}{1 - \kappa_{QG}} - \frac{\kappa_{12}\kappa_{GQ}(1 - \kappa_{21})}{1 - \kappa_{GQ}} \right)$$

$$f_2 = \frac{1}{1 - \kappa_{12}\kappa_{21}} \left(\frac{\kappa_{21}(1 - \kappa_{12})}{1 - \kappa_{QG}} - \frac{\kappa_{GQ}(1 - \kappa_{21})}{1 - \kappa_{GQ}} \right)$$

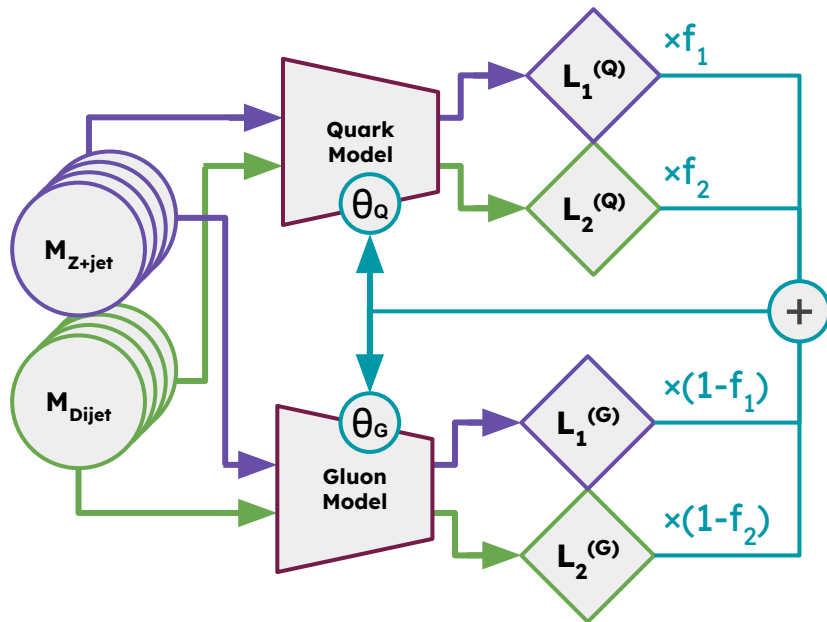
Backup: Reducibility correction



$$\kappa_{\text{QG}} = 1 - \frac{1 - \kappa_{12}}{f_1 - \kappa_{12}f_2}$$

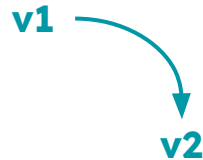
Backup: TopicFlow (v2)

- Joint training of the quark and gluon models gives convex loss:



$$p_Q(x) = \frac{p_{M_1}(x) - \kappa_{12} p_{M_2}(x)}{1 - \kappa_{12}}$$

$$p_G(x) = \frac{p_{M_2}(x) - \kappa_{21} p_{M_1}(x)}{1 - \kappa_{21}}$$

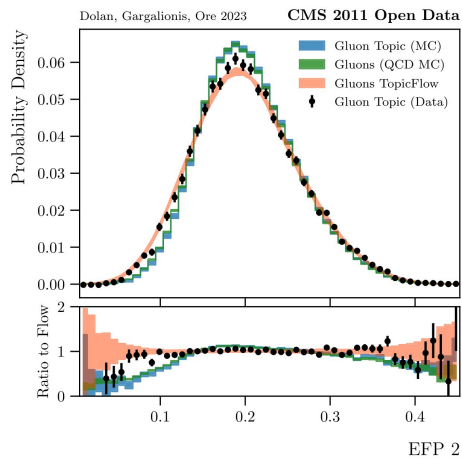
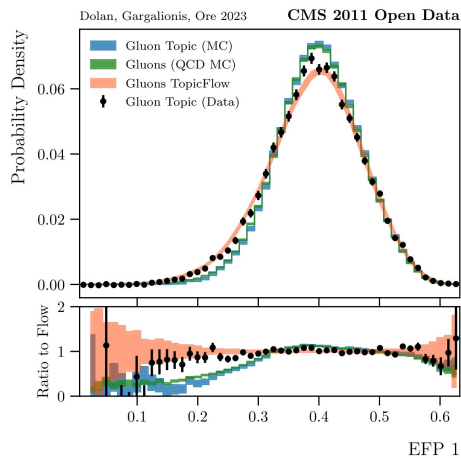
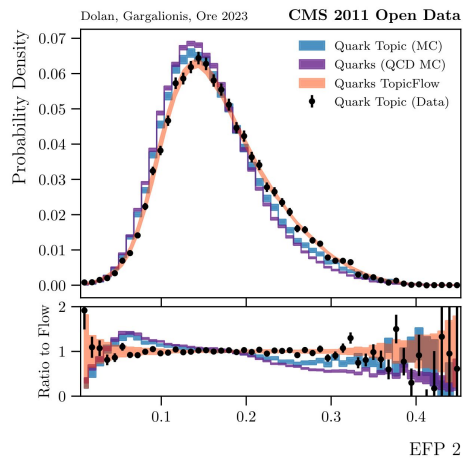
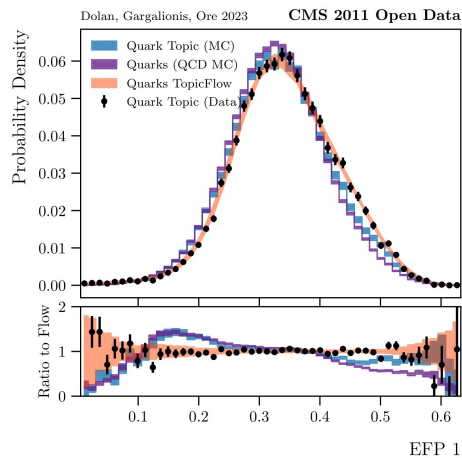


$$p_{M_1}(x) = f_1 p_Q(x) + (1 - f_1) p_G(x)$$

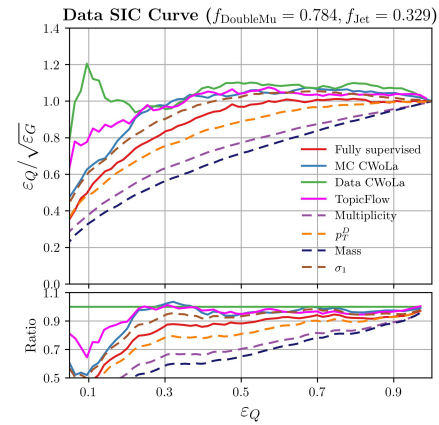
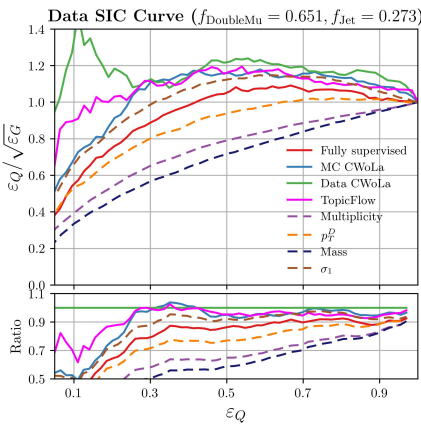
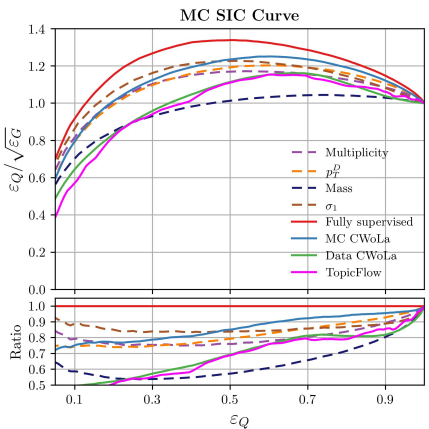
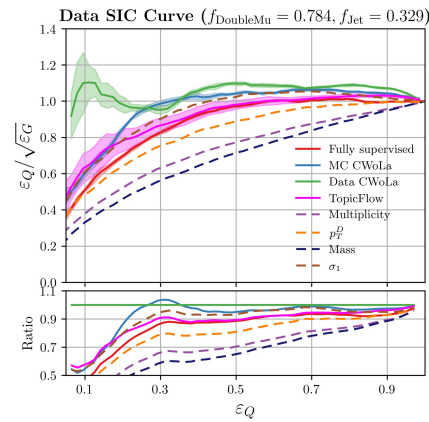
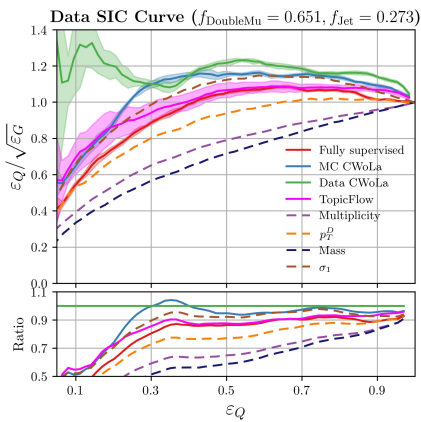
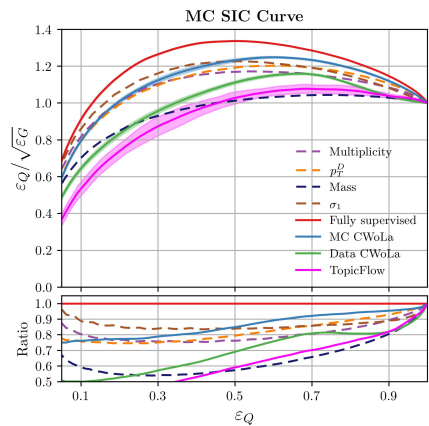
$$p_{M_2}(x) = f_2 p_Q(x) + (1 - f_2) p_G(x)$$

$$\begin{aligned} \text{Loss} &= \langle L_{\theta_Q, \theta_G}(x) \rangle_{x \sim p_{Z+\text{jet}}(x)} + \langle L_{\theta_Q, \theta_G}(x) \rangle_{x \sim p_{\text{dijet}}(x)} \\ &= f_1 \langle L_{\theta_Q}(x) \rangle_{x \sim p_{Z+\text{jet}}(x)} + (1 - f_1) \langle L_{\theta_G}(x) \rangle_{x \sim p_{Z+\text{jet}}(x)} \\ &\quad + f_2 \langle L_{\theta_Q}(x) \rangle_{x \sim p_{\text{dijet}}(x)} + (1 - f_2) \langle L_{\theta_G}(x) \rangle_{x \sim p_{\text{dijet}}(x)} \end{aligned}$$

Backup: TopicFlow samples



Backup: Ensemble tagging



Backup: Generative Classification

