

# MEM from the $ttH(bb)$ trenches

**Joosep Pata (ETH Zurich)**

Higgs Toppings

Benasque 2018

# Problem statement

- In a nutshell: given observed data  $\mathbf{y} \in \{\text{jets, leptons, ...}\}$ , exclude alternate hypotheses
- Search: **signal** ( $tt+H$ ,  $\mathbf{H}_1$ ) and **bkg** ( $tt+\text{jets}$ ,  $\mathbf{H}_0$ )

Neyman-Pearson lemma: given observable event  $\mathbf{y}$  and two hypotheses, the ratio

$$\lambda(\mathbf{y}) = \frac{L(\mathbf{y}|\mathbf{H}_1)}{L(\mathbf{y}|\mathbf{H}_0)}$$

is the **most powerful test statistic**, with  $L(\mathbf{y} | H)$  being the likelihood of data given observation

# Methods

We cannot directly compute the likelihood  $L(\mathbf{y} | H)$  given an observed event  $\mathbf{y}$ !

intractable integral:  $L(\mathbf{y} | H) \sim \int L(\mathbf{y} | z) p(z | H) dz$

The diagram shows the integral  $L(\mathbf{y} | H) \sim \int L(\mathbf{y} | z) p(z | H) dz$  with four labels and arrows: 'detector level' points to  $L(\mathbf{y} | z)$ ; 'theory amplitude' points to  $p(z | H)$ ; 'parton to detector transfer' points to the boundary between the two terms; and 'parton level' points to  $dz$ . The term  $L(\mathbf{y} | z)$  is highlighted in yellow and  $p(z | H)$  is highlighted in green.

We have sets of MC simulation **events with different H**:

$$\{\mathbf{y}_{1,0}\} \sim L(\mathbf{y} | H_{1,0})$$

- (I) Choose **clever observable**  $d=f(\mathbf{y})$ , estimate  $L(d(\mathbf{y}) | H)$  from simulation events: templates
- (II) Compute **approximate**  $L(\mathbf{y} | H)$  directly given theory ideas

# Clever observables

- Most analyses:  $d(\mathbf{y})$  is a parametric function (BDT, DNN), **numerically optimized** to discriminate signal from background based on MC samples
- **Amount of MC statistics needed?** Half a billion full-MC per experimental conditions is routine...
- **Which features exploited?** Study shaping of inv. masses, interpretability...
- **What if MC does not accurately represent data?**  
Cover with systematic uncertainties, but any bias?

*"Modelling studies can be seen which tend to overestimate certainty, pretending to produce crisp numbers precise to the third decimal digits even in situation of pervasive uncertainty or ignorance."*

jet-to-parton associations

parton level momenta

observed event

hypothesis

$$P(\mathbf{y}, \boldsymbol{\theta}) = \sum_{k=1}^{N_a} \int \frac{dx_1 dx_2}{2x_1 x_2 s} \int \prod_{i=1}^n \frac{d^3 p_i}{(2\pi)^3 2E_i} \times \delta^4(q_1 + q_2 - \sum_{i=1}^n p_i)$$

PDFs  $\times g(x_1)g(x_2)$

recoil  $\times \mathcal{R}(\tilde{\boldsymbol{\rho}}_T, \boldsymbol{\rho}_T)$

scattering amplitude  $\times |\mathcal{M}_{\boldsymbol{\theta}}(q_1, q_2, p_1, \dots, p_n)|^2$

transfer function  $\times W(\mathbf{y}, \mathbf{p})$ .

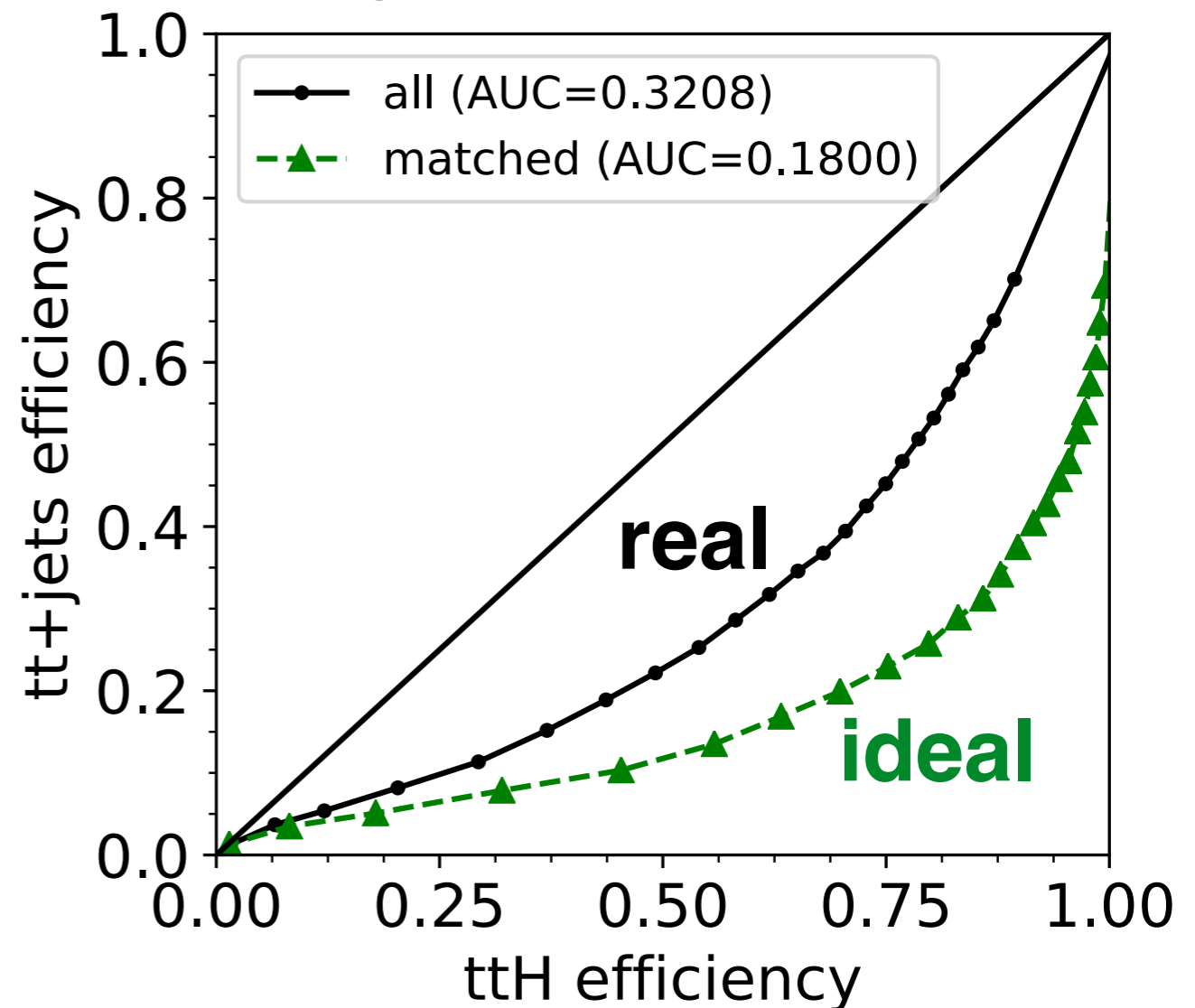
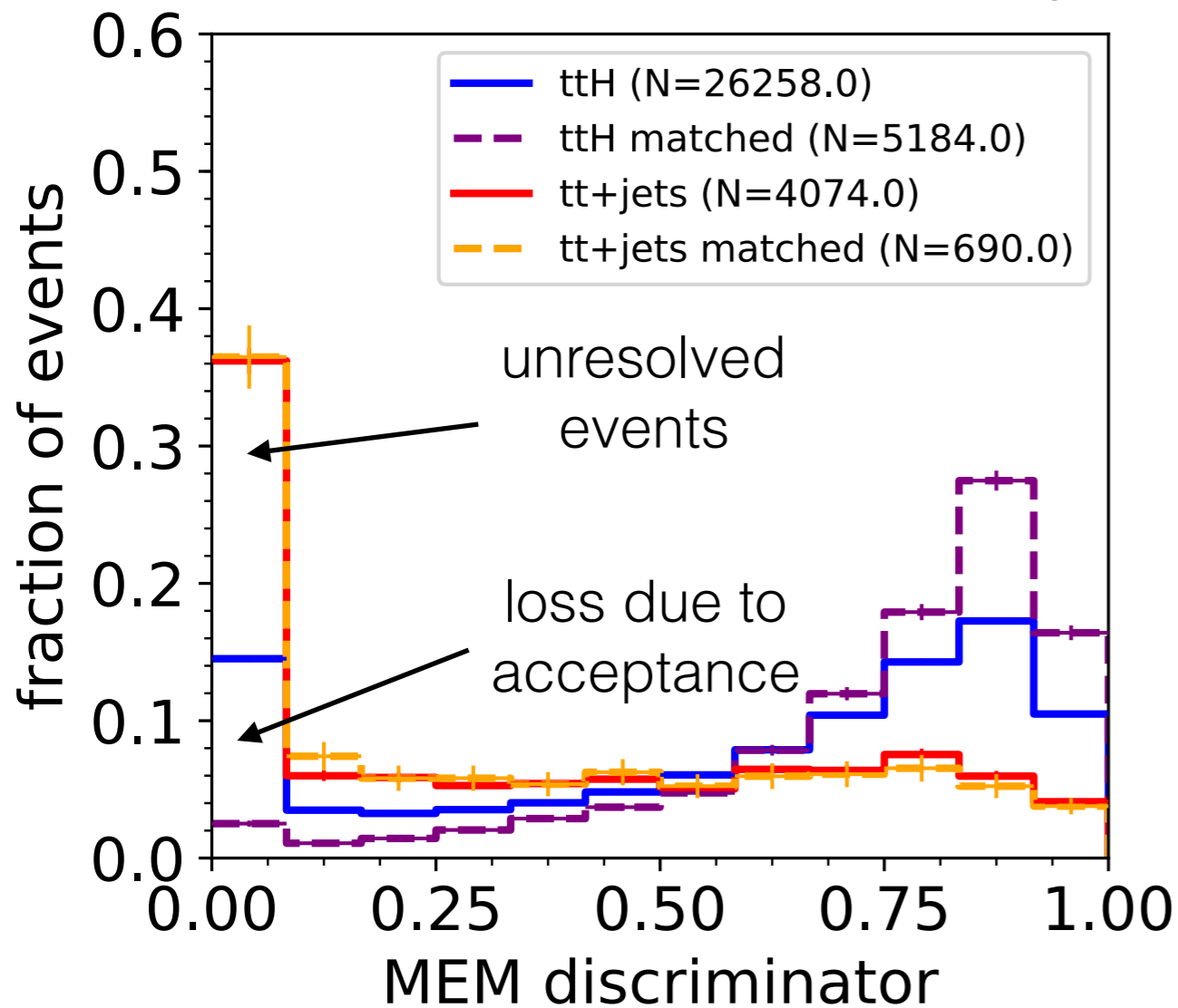
**Integrate numerically on the LHC computing grid.**

# Transfer functions

- Transition from parton-level event  $\mathbf{p}$  (a few 4-momenta) to detector-level event  $\mathbf{y}$  (~50-100 high-level quantities): showering, hadronization, detector effects, acceptance
- Assumptions, e.g. **Gaussian** quark-to-jet smearing
- Determined from **MC simulation by max likelihood fitting**
- **Need to hand-code** event-to-event transfer function, combination of possible jet-to-quark assignments
- In practice, integration becomes very expensive

# Reconstruction effects

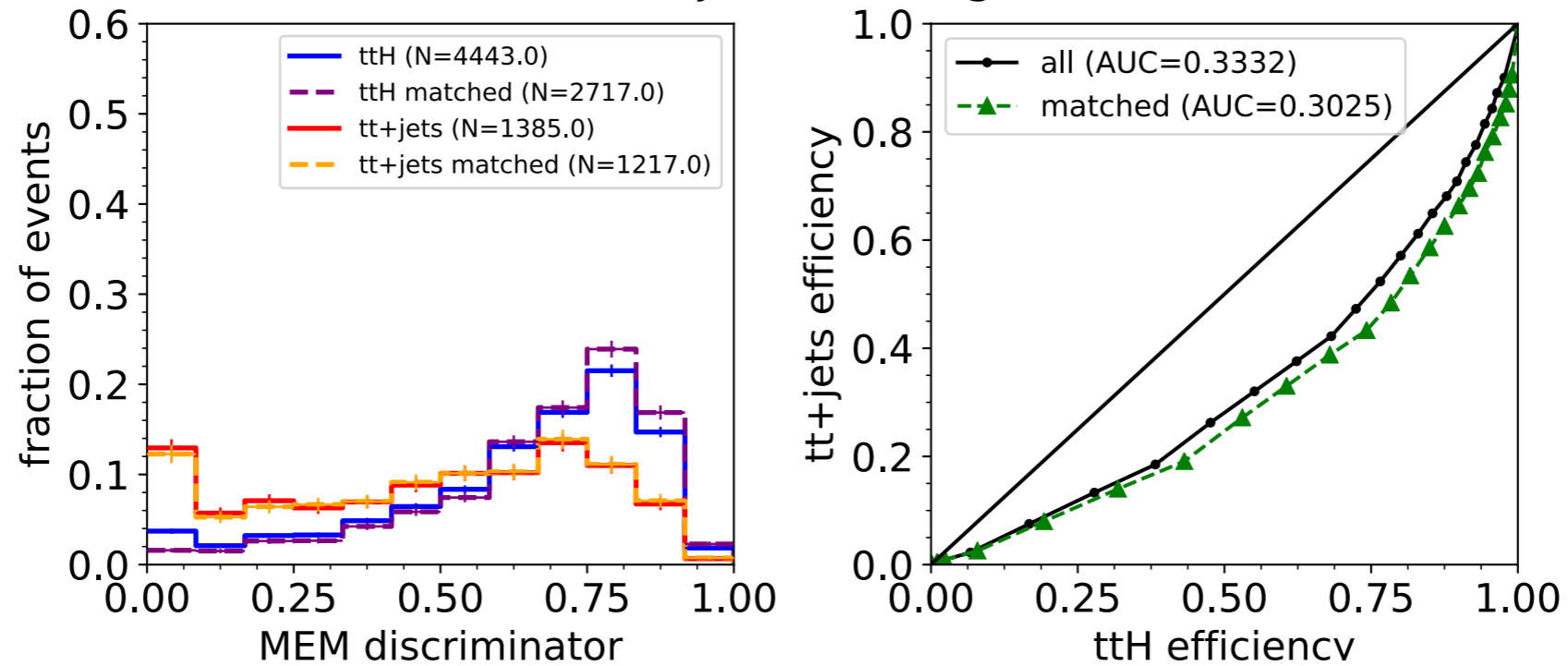
SL,  $\geq 6$  jets,  $\geq 4$  b-tags



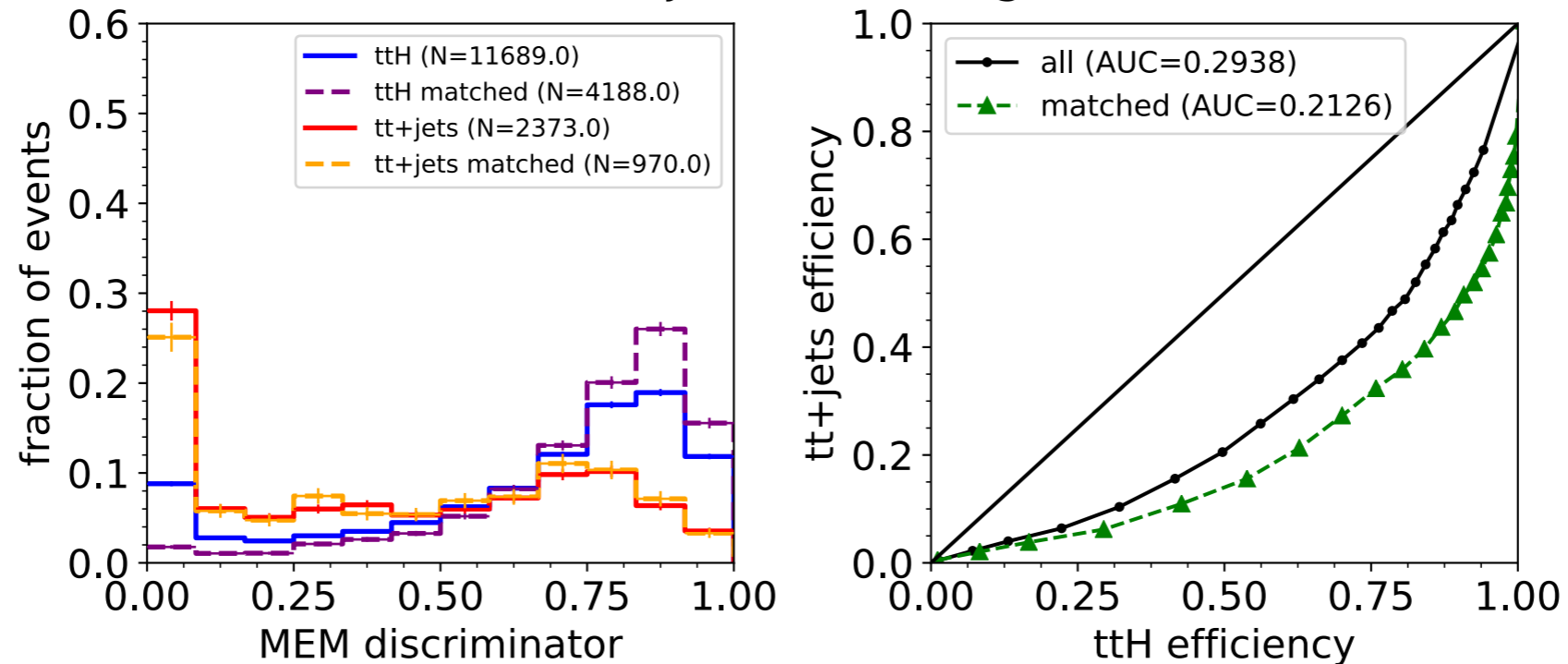
In ttH(bb), MEM performance is far from ideal due to detector effects & missed quarks!

# 4j, 5j events

SL, 4 jets, 4 b-tags



SL, 5 jets,  $\geq 4$  b-tags





# Questions for MEM

- **Which hypothesis?** ttH(bb) SL  $\geq 6j$ ,  $\geq 4b$  fully matched  $\sim 30\%$  of times, 30% miss quark from W
  - need to "degrade gracefully", but not implement MC by hand
- **Which assumptions?** Need to reduce integration space by assuming e.g. b-tagging, top tagging
- **More complex analyses?** Multi-hypothesis, multi-parameter EFT fits
- **NLO:** given computational cost, how to make use in experiment besides samples?

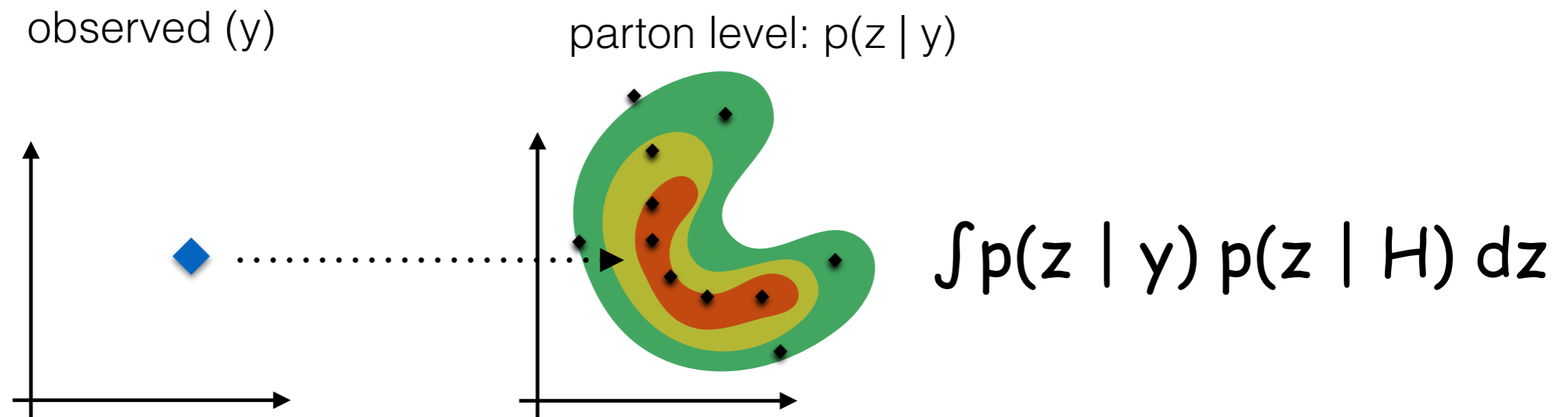
# MEM + ML combination

$$L(y | H) \sim \int L(y | z) p(z | H) dz$$

Use ML to approximate hadronization, showering, acceptance in  $L(y | z)$ .

Bonus: differentiable wrt. exp systematics, no reimplement of MC.

If possible to generate samples cheaply, can integrate scattering ME efficiently.



Or, replace per-event integral with regression over full MC:  
[Brehmer, Cranmer, Louppe, Pavez 2018]

# Bonus slides

# MEM hypotheses

interpretation	bottom quarks	light quarks	description
SL $2_W 2_h 2_t$	4	2	fully-reconstructed semileptonic
SL $1_W 2_h 2_t$	4	1	$(l\nu'_l b)_t (\not{q} q' \bar{b})_{\bar{t}} (b\bar{b})_h$
SL $0_W 2_h 2_t$	4	0	$(l\nu'_l b)_t (\not{q} \not{q}' \bar{b})_{\bar{t}} (b\bar{b})_h$
SL $2_W 2_h 2_t + 1g$	4	3	fully-reconstructed, additional ISR gluon
DL $2_h 2_t$	4	0	fully-reconstructed dileptonic

**Table 4.1:** The detailed event interpretations for semileptonic and dileptonic  $t\bar{t}H$  (signal) and  $t\bar{t} + b\bar{b}$  (background) events. In the semileptonic channel, we consider cases where up to 2 light quarks may be lost. The minimum number of jets required for a hypothesis is the sum of the number of quarks. For the SL  $1_W 2_h 2_t$  hypothesis the direction and energy of one of the light quarks is integrated out, denoted by  $\not{q}$ . For the fully-reconstructed semileptonic case with 7 jets, we also consider the ISR-modified interpretation.

# MEM permutations

interpretation	8+ jets	7 jets	6 jets	5 jets	4 jets
SL $2_W 2_h 2_t$	72/36	36	12	-	-
SL $1_W 2_h 2_t$	-	-	-	12	-
SL $0_W 2_h 2_t$	-	-	-	-	12
SL $2_W 2_h 2_t + 1g$	72/36	36	-	-	-
DL $2_h 2_t$	12	12	12	12	12

**Table 4.2:** The number of MEM combinations in associating quarks to jet for various MEM categories. In SL events with 7 or more jets, we choose up to 4 candidates based on invariant mass among the light jets for the W-boson reconstruction, in order to prevent a combinatorial explosion for events with a very high jet multiplicity. In DL events, we always choose exactly 4 candidates for the b quarks.

# CPU cost

method	time $t\bar{t}H$ (s)	time $t\bar{t} + b\bar{b}$ (s)	ROC AUC	$\epsilon_{\text{bkg}}$	total (h) / 1k
SL, $\geq 7\text{jet}, 2_W 2_h 2_t$	$45.8 \pm 18.9$	$69.4 \pm 26.1$	0.315	0.232	32.00
SL, $\geq 7\text{jet}, 2_W 2_h 2_t 1_g$	$71.7 \pm 27.1$	$471.7 \pm 50.6$	0.317	0.233	150.94
SL, $\geq 6\text{jet}, 2_W 2_h 2_t$	$30.2 \pm 21.0$	$45.4 \pm 30.9$	0.321	0.233	21.00
SL, $\geq 6\text{jet}, 1_W 2_h 2_t$	$64.8 \pm 22.7$	$101.1 \pm 33.0$	0.307	0.210	46.07
SL, $\geq 6\text{jet}, 0_W 2_h 2_t$	$83.9 \pm 20.4$	$136.3 \pm 28.9$	0.294	0.218	61.16
SL, 5jet, $1_W 2_h 2_t$	$25.4 \pm 7.1$	$39.9 \pm 9.6$	0.293	0.198	18.13
SL, 5jet, $0_W 2_h 2_t$	$84.7 \pm 20.3$	$136.9 \pm 28.9$	0.291	0.217	61.56
SL, 4jet, $0_W 2_h 2_t$	$84.3 \pm 20.7$	$136.0 \pm 29.2$	0.333	0.275	61.21
DL, $\geq 4\text{jet}, 0_W 2_h 2_t$	$55.7 \pm 13.7$	$90.4 \pm 19.3$	0.223	0.124	40.58

**Table 4.3:** The CPU budget and separation power of the MEM in the SL channel using various event interpretations. We show the time required to evaluate the signal and background hypotheses, the receiver operating characteristic (ROC) area under curve (AUC), the efficiency of background at a signal efficiency of 50% ( $\epsilon_{\text{bkg}}$ ) and the total time required to compute the MEM discriminator for 1000 events.