# Neural Inference:

### what does it mean?
### why should I care?
### why can't I have it right now?

Chris Pollard, Warwick

before we can talk about neural inference,

we have to believe machine-learning is doing what it claims!

# Multilayer Feedforward Networks are Universal Approximators

KURT HORNIK

Technische Universität Wien

MAXWELL STINCHCOMBE AND HALBERT WHITE

University of California, San Diego

Abstract—This paper rigorously establishes that standard multilayer feedforward networks with as few as one hidden layer using arbitrary squashing functions are capable of approximating any Borel measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many hidden units are available. In this sense, multilayer feedforward networks are a class of universal approximators.

NNs can approximate any continuous function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$

**to *arbitrary precision,***

given a finite number of free parameters
(+ enough training data and time).

3

neural inference:

extracting parameters of interest from observations
using NN approximations of a likelihood(-ratio) or posterior.

neural inference:

extracting parameters of interest from observations
using NN approximations of a likelihood(-ratio) or posterior.

how can we achieve this?

mostly it's a matter of constructing (and minimizing)
*the correct loss function*.

with a *simulator s* that,

given a parameter $\phi$,

produces observables $x$

$$s : \phi \to p(x) \equiv p(x \,|\, \phi),$$

with a *simulator* $s$ that,

given a parameter $\phi$,

produces observables $x$

$$s : \phi \rightarrow p(x) \equiv p(x \,|\, \phi),$$

if we have many simulated samples $(\phi, x)$,

then we know how to learn the log-likelihood ratio

$$D_{\phi'}^{\phi}(x) \equiv \log \frac{p(x \,|\, \phi)}{p(x \,|\, \phi')}.$$

we know how to learn the log-likelihood ratio

$$D^{\phi}_{\phi'}(x) \equiv \log \frac{p(x \mid \phi)}{p(x \mid \phi')}.$$

imagine the case that $\phi$ only has two possible values: $A, B$.

$$D^{A}_{B}(x) \equiv \log \frac{p(x \mid A)}{p(x \mid B)}$$

is simply what a *discriminator* between $A$ and $B$ learns to emit.

we ... ow to learn the log-likelihood ratio

$$D_{\phi'}^{\phi}(x) \equiv \log \frac{p(x\,|\,\phi)}{p(x\,|\,\phi')}.$$

imagine the case that $\phi$ only has two possible values: $A, B$.

$$D_B^A(x) \equiv \log \frac{p(x\,|\,A)}{p(x\,|\,B)}$$

is simply what a *discriminator* between $A$ and $B$ learns to emit.

given a simulator $p(x \,|\, \phi)$, we also know, how to learn

$$p(\phi \,|\, x) \propto p(x \,|\, \phi) \, p(\phi)$$

for some family of posteriors $p_\eta$ parameterized by $\eta$:

given a simulator $p(x\,|\,\phi)$, we also know, how to learn

$$p(\phi\,|\,x) \propto p(x\,|\,\phi)\,p(\phi)$$

for some family of posteriors $p_\eta$ parameterized by $\eta$:

$p(\phi)$ is the distribution over $\phi$ used in training,

and we learn the function $f : x \to \eta$ by minimizing the loss

$$L(\phi, x) = -\log p_{\eta=f(x)}(\phi).$$

giv... ...tor $p(x|\phi)$, we also know, how to learn

$$p(\phi|x) \propto p(x|\phi)\, p(\phi)$$

for some family of posteriors $p_\eta$ parameterized by $\eta$:

$p(\phi)$ is the distribution over $\phi$ used in training,

and we learn the function $f : x \rightarrow \eta$ by minimizing the loss

$$L(\phi, x) = -\log p_{\eta=f(x)}(\phi).$$

summary:

*if* we believe NNs work the way they should…

then given a simulator $p(x \,|\, \phi)$

we can learn an *arbitrarily precise approximation* of

$$\log \frac{p(x \,|\, \phi)}{p(x \,|\, \phi')} \quad \text{or} \quad p(\phi \,|\, x).$$

summary:

*if* we believe NNs work the way they should…

then given a simulator $p(x|\phi)$

arbitrarily complex!

we can learn an *arbitrarily precise approximation* of

$$\log \frac{p(x|\phi)}{p(x|\phi')} \quad \text{or} \quad p(\phi|x).$$

summary:

**if** we believe NNs work the way they should…

then given a simulator $p(x|\phi)$

arbitrarily complex!

we can learn an *arbitrarily precise approximation* of

$$\log \frac{p(x|\phi)}{p(x|\phi')} \quad \text{or} \quad p(\phi|x).$$

leads to *best possible* constraints

summary:

*if* we believe NNs work the way they should…

then given a simulator $p(x|\phi)$

arbitrarily complex!

we can learn an *arbitrarily precise approximation* of

the entire analysis is "basically done."

$$\log\frac{p(x|\phi)}{p(x|\phi')} \quad \text{or} \quad p(\phi|x).$$

leads to *best possible* constraints

summary:

**if** we believe NNs work the way they should…

then given a simulator $p(x|\phi)$ | arbitrarily complex! |

we can learn an *arbitrarily precise approximation* of

the entire analysis is "basically done."

$$\log \frac{p(x|\phi)}{p(x|\phi')} \quad \text{or} \quad p(\phi|x).$$

leads to *best possible* constraints

17

too good to be true?

too good to be true?

at least three* open "issues"
with neural simulation-based inference (nSBI):

*convergence of function approximations*

*global nuisance parameters*

*model mis-specification*

# *convergence of function approximations*



adobe stock

# *convergence of function approximations*

the likelihood-ratio over an entire dataset
may require very, very good approximations
of the per-event likelihood-ratio to behave properly.

$$L(x \mid \phi) = \prod_i L(x_i \mid \phi)$$

(for independent observations)

# *convergence of function approximations*

the likelihood-ratio over an entire dataset
may require very, very good approximations
of the per-event likelihood-ratio to behave properly.

$$L(x\,|\,\phi) = \prod_i L(x_i\,|\,\phi)$$

(for independent observations)

this may be *difficult*,
but at least it can be *verified* on toy datasets.

*global
nuisance parameters*

# *global nuisance parameters*

in the presence of global nuisance parameters,
per-event likelihoods do not easily compose.

we need to be able to profile against
or marginalize over global nuisance parameters.

# *global nuisance parameters*

two main approaches:

1) learn the likelihood/posterior *parameterized* in the nuisances and profile/marginalize after the fact, or

2) perform dataset-wide (ensemble) learning.

# *global nuisance parameters*

**works but can be painful**

two main approaches:

1) learn the likelihood/posterior *parameterized* in the nuisances and profile/marginalize after the fact, or

2) perform dataset-wide (ensemble) learning.

arXiv:2412.01548



arXiv:2412.01600

ATLAS recently published a measurement of off-shell Higgs production using this approach:

→ it does work!

→ very high CPU costs

→ currently limited to relatively small datasets.

# *global nuisance parameters*

two main approaches:

**works but can be painful**

1) learn the likelihood/posterior *parameterized* in the nuisances and profile/marginalize after the fact, or

2) perform dataset-wide (ensemble) learning.

**has some very nice properites**

$x_i$ → $s$ → $y_i$

$x_j$ → $s$ → $y_j$

$x_k$ → $s$ → $y_k$

$\Sigma$

$\rho$ → $p(\phi \,|\, \{x\})$

per-event network

permutation invariant aggregator

dataset-wide inference network

inference over an *ensemble*, $\{x\}$,
via the *deep set* architecture

derive a vector summary observable per event

$x_i$ → [ $s$ ] → $y_i$

$x_j$ → [ $s$ ] → $y_j$

$x_k$ → [ $s$ ] → $y_k$

$\Sigma$

[ $\rho$ ] → $p(\phi \,|\, \{x\})$

per-event
network

permutation
invariant
aggregator

dataset-wide
inference
network

sum over these per-event summaries

$x_i$ → $s$ → $y_i$

$x_j$ → $s$ → $y_j$

$x_k$ → $s$ → $y_k$

$\Sigma$

$\rho$ → $p(\phi \,|\, \{x\})$

per-event network

permutation invariant aggregator

dataset-wide inference network

perform variational inference on this dataset-wide summary

$x_i$ → $s$ → $y_i$

$x_j$ → $s$ → $y_j$

$x_k$ → $s$ → $y_k$

$\sum$

$\rho$ → $p(\phi \,|\, \{x\})$

per-event network

permutation invariant aggregator

dataset-wide inference network

very simple case: inferring diagonal components of $\vec{\sigma}$ from observations

$$x \sim \mathcal{N}(\vec{\mu}, \mathrm{diag}(\vec{\sigma}^2))$$

very simple case: inferring  diagonal
components of $\vec{\sigma}$ from observations

$$x \sim \mathcal{N}(\vec{\mu}, \text{diag}(\vec{\sigma}^2))$$

very simple case: inferring diagonal components of $\vec{\sigma}$ from observations

$$x \sim \mathcal{N}(\vec{\mu}, \mathrm{diag}(\vec{\sigma}^2))$$

particle-physics "inspired" example:

goal: extract signal fraction $\theta$,

without knowing *a priori* the signal location, $\theta_\nu$.
($\theta_\nu$ is a global NP.)

plus a wide background contribution.

particle-physics "inspired" example:

goal: extract signal fraction $\theta$,

without knowing *a priori* the signal location, $\theta_\nu$.
($\theta_\nu$ is a global NP.)

plus a wide background contribution.

Posterior standard deviation

$\theta_\nu = \theta_{\nu, \text{nom}}$

Deep set
MCMC on $x$
MCMC on $s_{\text{nom}}$
MCMC on $s_{\text{marg}}$

$N_0 = 100$
$\theta_{\text{true}} = 0.2$

med[$\sigma_\theta$]

True signal position $\theta_{\nu, \text{true}}$

as we sweep $\theta_\nu$,

the **deep-set posterior**

agrees with

**brute-force (MCMC)**

for posterior estimation on $\{x\}$.

**Posterior standard deviation**

discriminators learned per-event on $x$

**at the nominal signal position**

**and marginalized over $\theta_\nu$**

*do not preserve enough information to build the correct posterior $\forall\, \theta_\nu$.*

Posterior standard deviation

$\theta_\nu = \theta_{\nu, \text{nom}}$

Deep set
MCMC on $x$
MCMC on $s_{\text{nom}}$
MCMC on $s_{\text{marg}}$

$N_0 = 100$
$\theta_{\text{true}} = 0.2$

med[$\sigma_\theta$]

True signal position $\theta_{\nu, \text{true}}$

the deep set is
*many orders of
magnitude faster to run*

than MCMC

or a "conventional"
profiled likelihood fit.

$x_i$ — $s$ → $y_i$

$x_j$ — $s$ → $y_j$

$x_k$ — $s$ → $y_k$

$\Sigma$ → $\rho$ →

coming back to
our network architecture

the *per-event embedding*

$$y \equiv s(x)$$

is information-preserving w.r.t. the posterior or likelihood over $\theta$,

*even in the presence of global NPs.*

updating the posterior with new data is simple:

$$y_{\{x\}+x_0} = y_{\{x\}} + y_{x_0}$$

then feed this to $\rho$.

$\rightarrow$ this is extremely fast: can run full inference ~in real time.

updating the posterior data is simple:

$$y_{\{x\}+x_0} = y_{\{x\}} + y_{x_0}$$

then feed this to $\rho$.

$\rightarrow$ this is extremely fast: can run full inference ~in real time.

updating the posterior with new data is simple:

$$y_{\{x\}+x_0} = y_{\{x\}} + y_{x_0}$$

then feed this to $\rho$.

$\rightarrow$ this is extremely fast: can run full inference ~in real time.

**ongoing work: scaling up to datasets with millions of events +**

45

*model mis-specification*

https://animegenius.live3d.io/features/dragon-ai-art-generator

# *model mis-specification*

summary:

*if* we believe NNs work the way they should…

then given a simulator $p(x|\phi)$ — arbitrarily complex!

we can learn an *arbitrarily precise approximation* of

the entire analysis is basically "done."

$\log \dfrac{p(x|\phi)}{p(x|\hat{\phi})}$ or $p(\phi|x)$.

leads to *best* possible constraints

we claimed for

"arbitrarily complex $x$ and $\phi$"

given a simulator $p(x|\phi)$

we can perform "correct" inference.

# *model mis-specification*

is the simulation
*arbitrarily trustworthy*?

summary:

*if* we believe NNs work the way they should…

then given a simulator $p(x \mid \phi)$ **arbitrarily complex!**

we can learn an *arbitrarily precise approximation* of

**the entire analysis is basically "done."**
$\log \dfrac{p(x \mid \phi)}{p(x \mid \hat{\phi})}$ or $p(\phi \mid x)$.
**leads to *best* possible constraints**

# *model mis-specification*

summary:

*if* we believe NNs work the way they should...

then given a simulator $p(x|\phi)$   **arbitrarily complex!**

we can learn an *arbitrarily precise approximation* of

**the entire analysis is basically "done."**   $\log \dfrac{p(x|\phi)}{p(x|\hat{\phi})}$   or   $p(\phi|x)$.   **leads to *best* possible constraints**

is the simulation *arbitrarily trustworthy*?

of course the answer is "no"!

how do we validate/calibrate very complex observables $x$ that enable strong, ***correct*** constraints?

# *model mis-specification*

summary:

**if** we believe NNs work the way they should…

then given a simulator $p(x|\phi)$    **arbitrarily complex!**

we can learn an *arbitrarily precise approximation* of

the entire analysis is
$\log \frac{p(x|\phi)}{p(x|\hat{\phi})}$ or $p(\phi|x)$.

**leads to *best* possible**

in general a very difficult problem, but we are working on it…

is the simulation *arbitrarily trustworthy*?

of course the answer is "no"!

how do we validate/calibrate

very complex observables $x$

that enable strong,
***correct*** constraints?

**we have a recipe for reweighting a density $p(x)$ (sim) to $q(x)$ (data):**

⊙ train a discriminator $D_p^q(x)$ and weight each sample $x$ by $\exp D_p^q(x)$.

⊙ allows weight-based calibration for complex observables.

**we have a recipe for reweighting a density $p(x)$ (sim) to $q(x)$ (data):**

◉ train a discriminator $D_p^q(x)$ and weight each sample $x$ by $\exp D_p^q(x)$.

◉ allows weight-based calibration for complex observables.

**in many instances it is preferable to *move* events to perform a calibration:**

$$T : \vec{x} \rightarrow \vec{x} \implies T_{\#}p \approx q$$

◉ for any $p$ and $q$, at least one $T$ exists, and it is unitary.

◉ we usually want to change the simulation *as little as possible*:

✦ we need to find the unique $\hat{T}$ that minimally (or "optimally") morphs $p$ into $q$.

**we have been implicitly using optimal transport (OT) for decades in HEP:**

◉ correct simulated scale of a gaussian density ↔ move the mean of simulated distribution.

**neural OT generalizes this process.**

◉ for euclidean spaces, the OT map is the gradient of some convex potential:

$$\hat{T}\,\vec{x} \equiv \vec{\nabla}\,\phi(\vec{x}).$$

◉ we have recently managed to learn $\phi_z$ and therefore $\hat{T}_z$ in concrete use-cases:

✦ the OT map is *conditional* on $z$.

**simulation**

**data**

$\vec{x}$    $\vec{x}'$

**we have been implicitly using optimal transport (OT) fo** ... **IEP:**

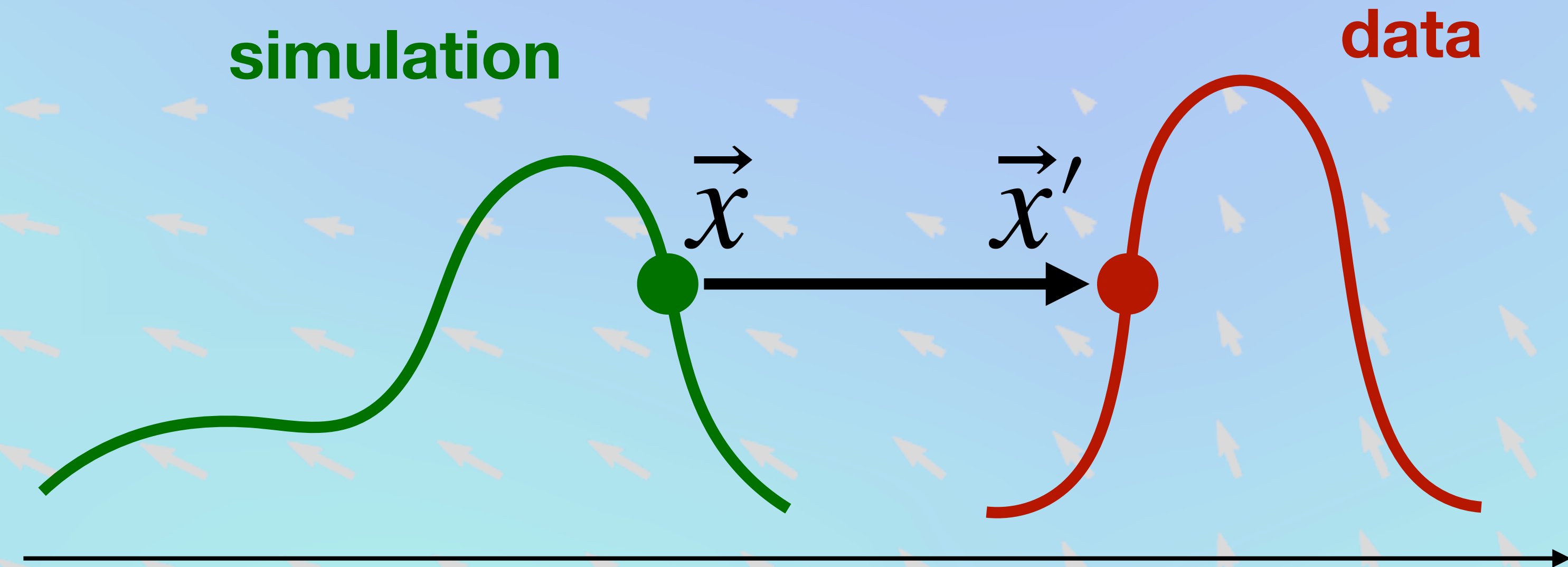◉ correct simulated scale of a gaussian density ↔ move the mean ... tribution.

**neural OT generalizes this process.**

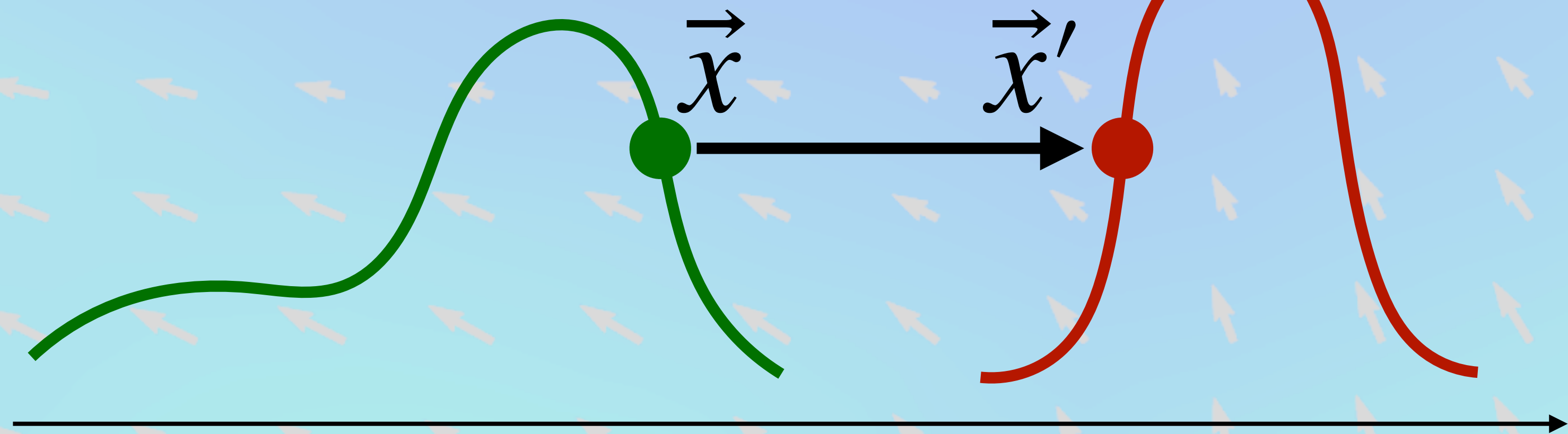◉ for euclidean spaces, the OT map is the gradient of some convex potential:

$$\hat{T}\,\vec{x} \equiv \vec{\nabla}\,\phi(\vec{x}).$$

◉ we have recently managed to learn $\phi_z$ and therefore $\hat{T}_z$ in concrete use-cases:

✦ the OT map is *conditional* on $z$.

**simulation**

**data**

$\vec{x}$      $\vec{x}'$

# a concrete example: jet flavor-tagging in ATLAS

# jet flavor-tagging is a classification problem:

◉ ATLAS's classifiers emit the probability of a jet to contain a $b$-hadron, $c$-hadron, or neither $(p_b, p_c, p_u)$.

◉ modern algorithms are very complex:

　✦ charged-particle tracks as inputs.

　✦ incredible separating power, but clear mismodeling

◉ until now, we had no direct calibration for these probabilities.

　✦ to do so, we set $q_i \equiv \mathrm{logit}\, p_i$, treat $\vec{q}$ as euclidean, and calibrate via OT.



56

notation:

- $q_i \equiv \mathrm{logit}\, p_i$ : flav. class. scores
- $p_{\mathrm{sim}}(\vec{q}\,|\,p_T) \equiv p_{\mathrm{sim}}(\vec{q}\,|\,p_T)\, p_{\mathrm{data}}(p_T)$
- $\hat{T}_{\#} \equiv p_T$-dependent OT map

**result:**

◉ we obtain the full 3D OT maps in
$\vec{q}$ space s.t. $\hat{T}_{\#} p_{\mathrm{sim}} \approx p_{\mathrm{data}}$,

✦ derived as a function of jet $p_T$.

◉ here we show a 2D slice for
$q_b \times q_c$ at fixed $q_u \times p_T$.



ATLAS Preliminary
$p_T = 64.0$ GeV, logit $p_u$=0.0, $\sqrt{s} = 13$ TeV, 140 fb$^{-1}$

notation:

- $q_i \equiv \mathrm{logit}\, p_i$ : flav. class. scores
- $p_{\mathrm{sim}}(\vec{q}\,|\,p_T) \equiv p_{\mathrm{sim}}(\vec{q}\,|\,p_T)\, p_{\mathrm{data}}(p_T)$
- $\hat{T}_{\#} \equiv p_T$-dependent OT map

**the technology works!**

◉ before calibration, poor modeling of many
quantities…

    ✦ including the $b$-tagging discriminator

$$D_b \equiv \log \frac{p_b}{f_c p_c + (1 - f_c) p_u}$$

◉ after calibration, very good agreement even for this
non-trivial quantity:

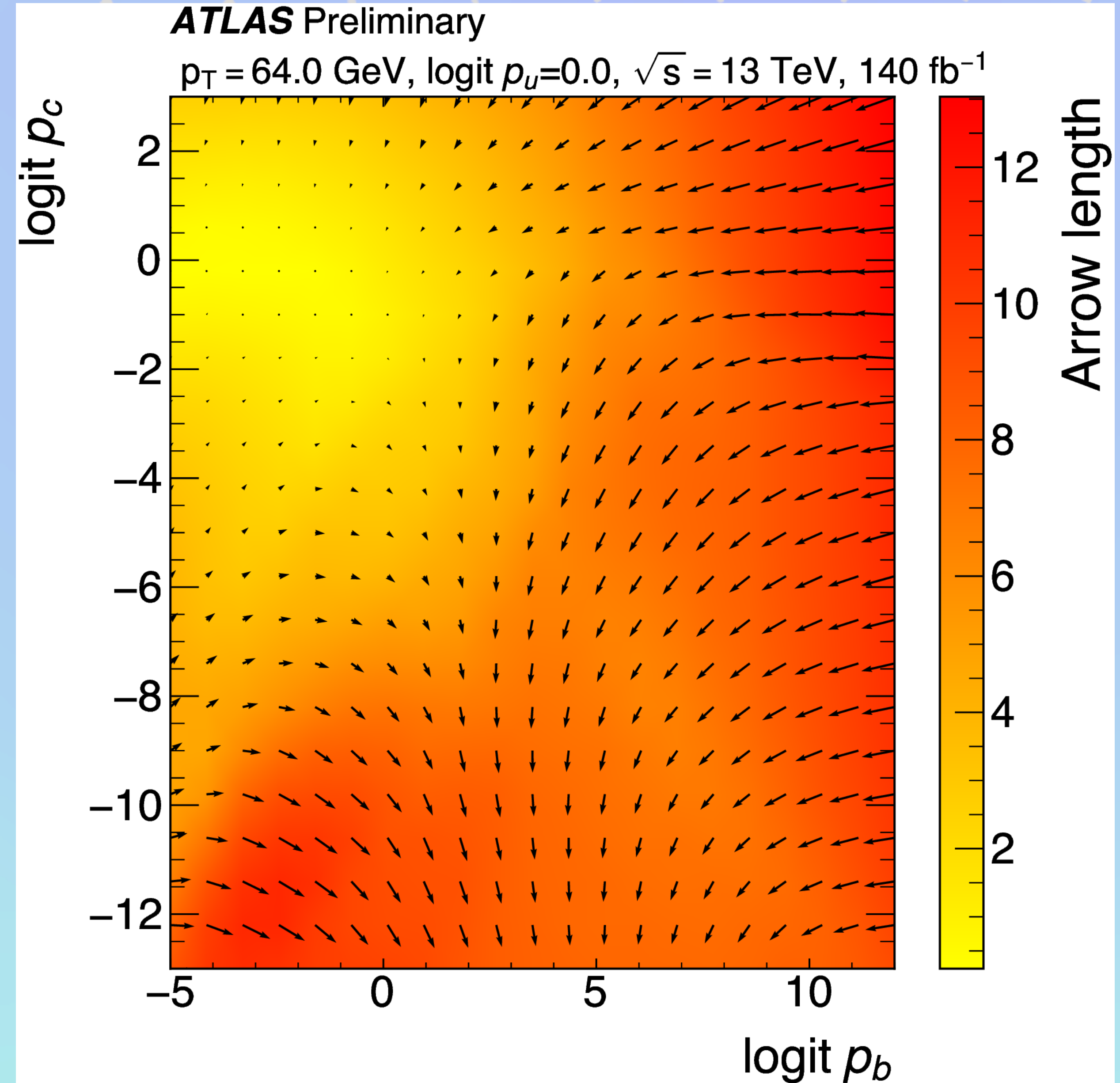    ✦ higher-order correlations between $p_i$ are
properly corrected.

notation:

- $q_i \equiv \text{logit}\, p_i$ : flav. class. scores
- $p_{\text{sim}}(\vec{q}\,|\,p_T) \equiv p_{\text{sim}}(\vec{q}\,|\,p_T)\, p_{\text{data}}(p_T)$
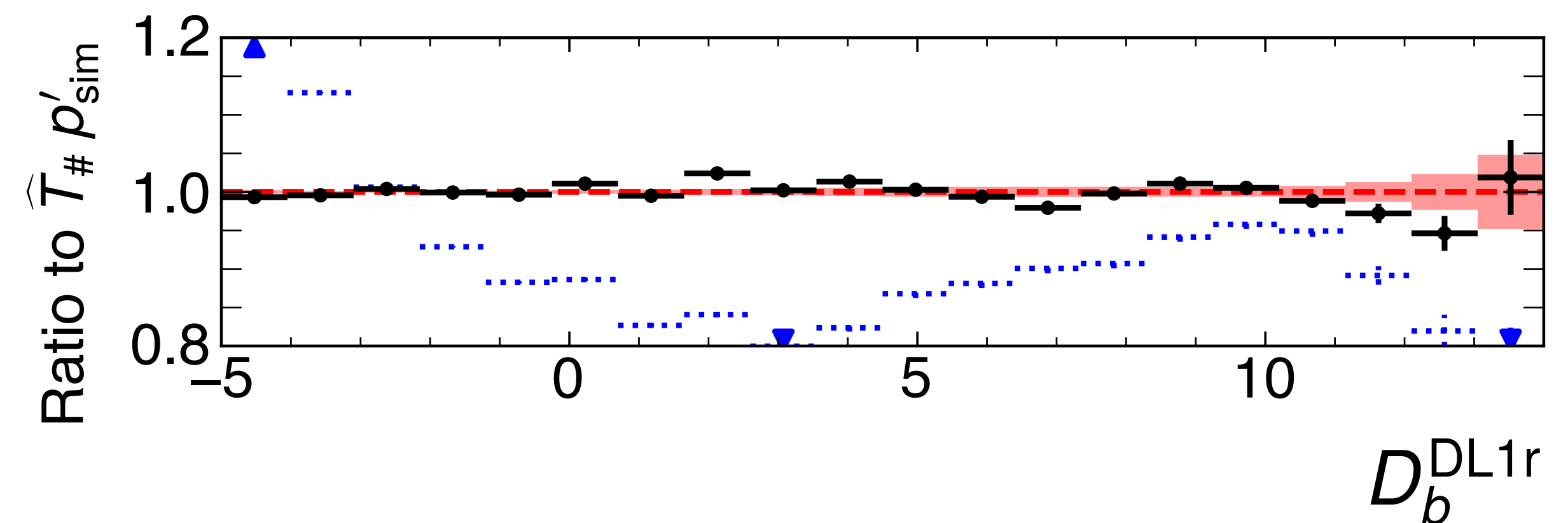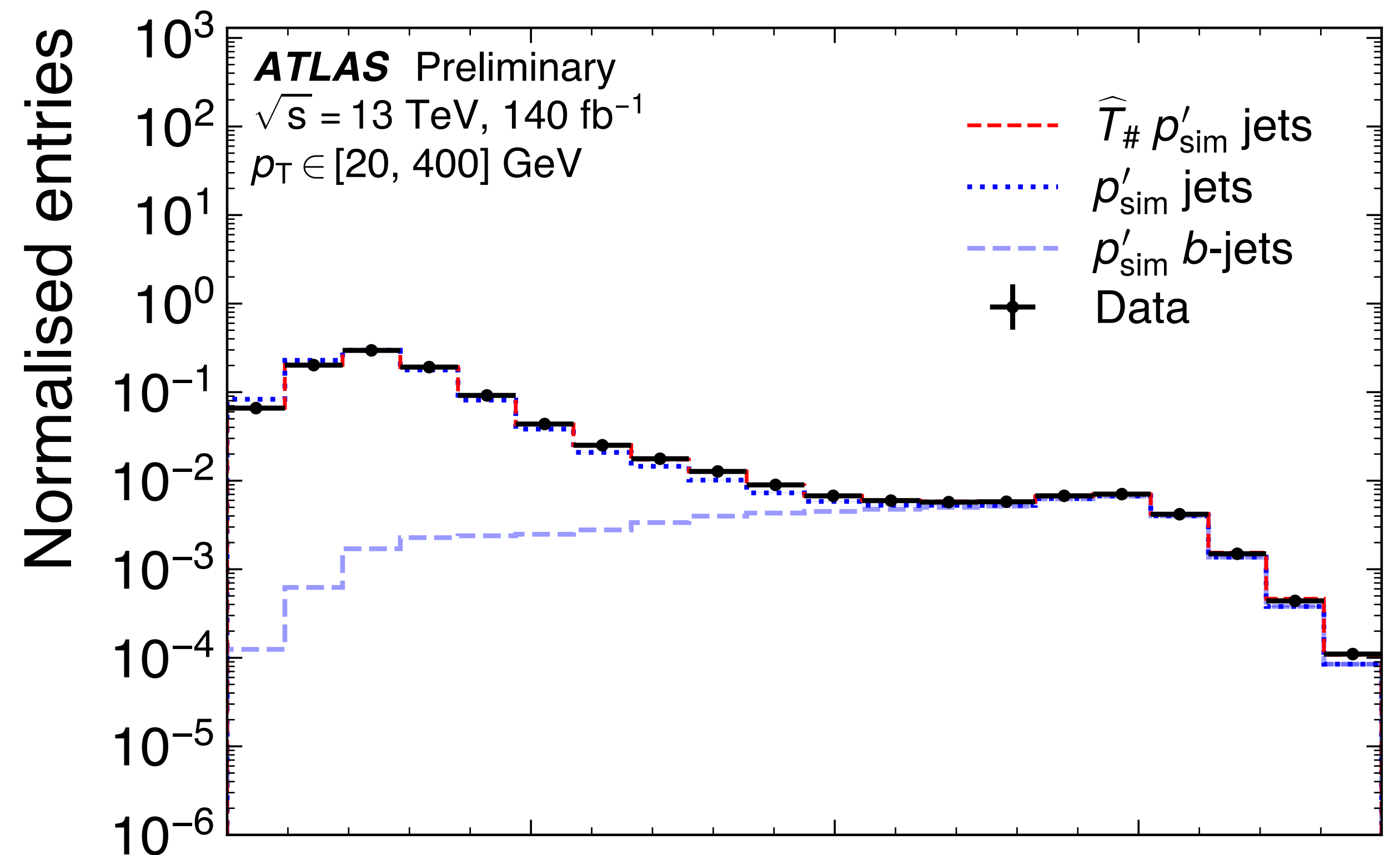- $\hat{T}_\# \equiv p_T$-dependent OT map

**the technology works!**

◉ before calibration, poor modeling of many quantities…

✦ including the $b$-tagging discriminator

$$D_b \equiv \log \frac{p_b}{f_c p_c + (1 - f_c)p_u}$$

◉ after calibration, very good agreement even for this non-trivial quantity:

✦ higher-order correlations between $p_i$ are properly corrected.



ATLAS Prelim
$\sqrt{s}$ = 13 TeV, 140 fb
$p_T \in [20, 400]$ GeV

$\hat{T}_\# \, p'_{\text{sim}}$ jets

Normalised entries

Ratio to $\hat{T}_\# \, p'_{\text{sim}}$

$D_b^{\text{DL1r}}$

# the future (?): scaling this up

# ongoing work: scaling it up

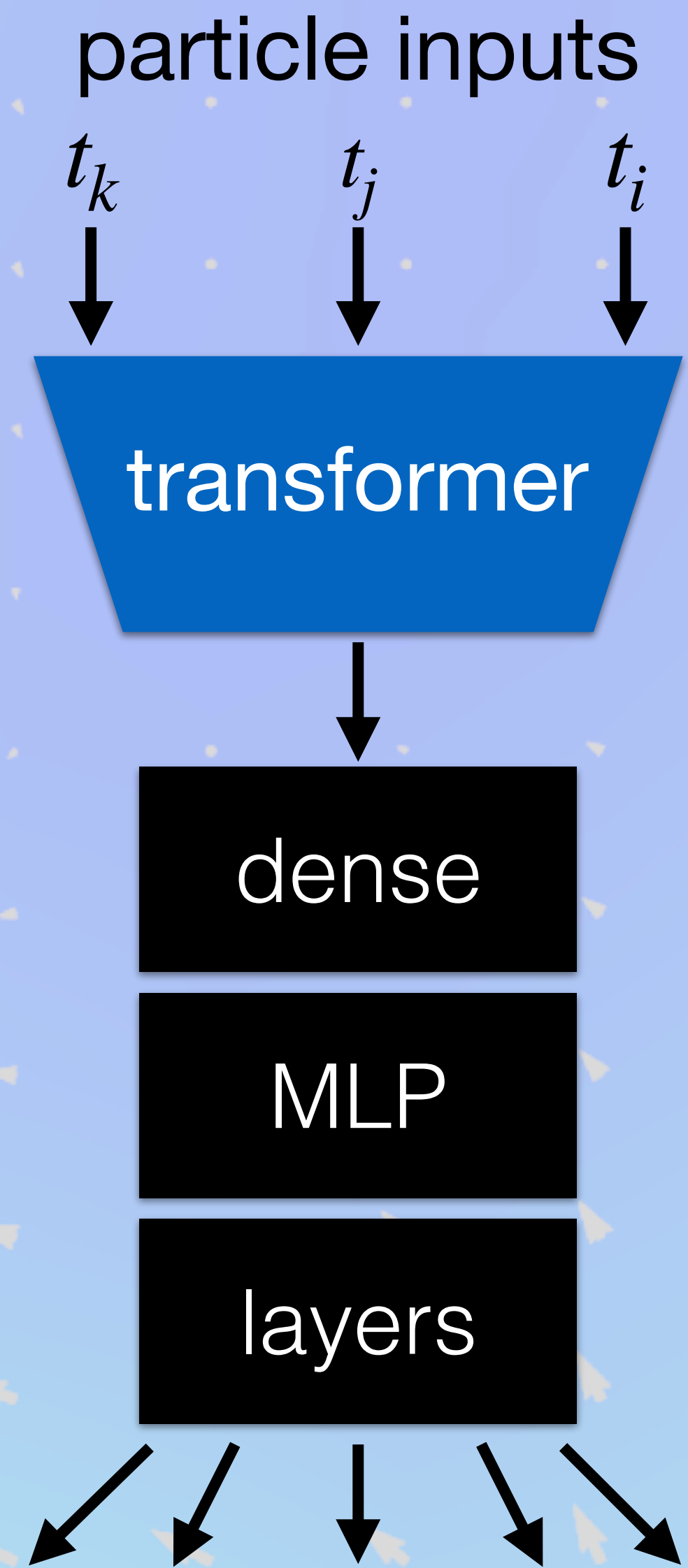◉ preliminary results show that this scales nicely to *very* high dimensions:

◉ start with a JetClass-like classifier that discriminates between 10 distinct types of large-radius jets:

✦ $H \to bb, H \to cc, t \to bqq'$, inclusive ("QCD") jets, etc

particle inputs
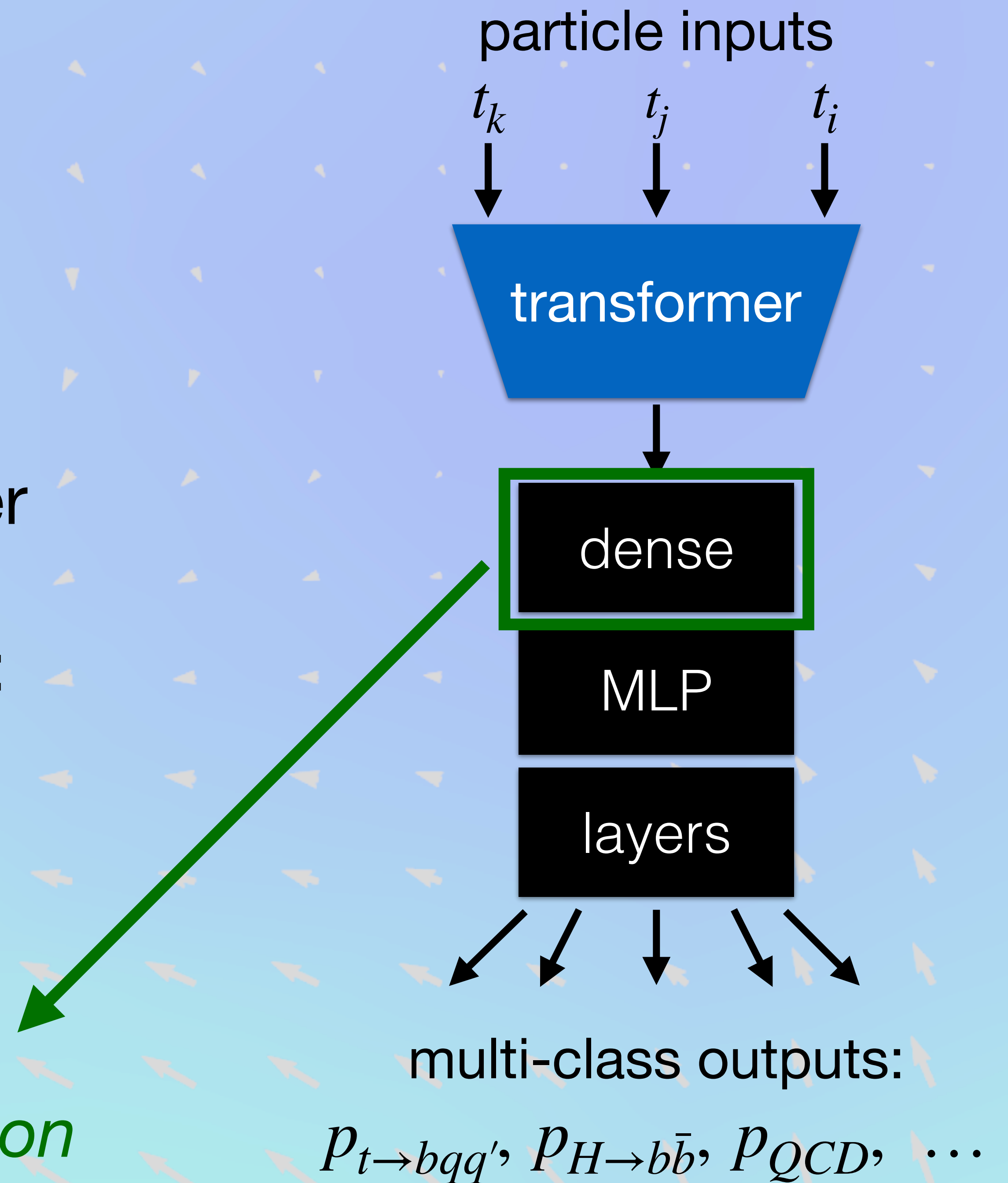
$t_k$     $t_j$     $t_i$

transformer

dense

MLP

layers

multi-class outputs:

$p_{t \to bqq'}, \ p_{H \to b\bar{b}}, \ p_{QCD}, \ \dots$

# ongoing work: scaling it up

◉ preliminary results show that this scales nicely to *very* high dimensions:

◉ start with a JetClass-like classifier that discriminates between 10 distinct types of large-radius jets:

✦ $H \to bb, H \to cc, t \to bqq'$, inclusive ("QCD") jets, etc

◉ *calibrate the internal 128-dim representation of the jet information*

particle inputs

$t_k$     $t_j$     $t_i$

transformer

dense

MLP

layers

multi-class outputs:

$p_{t \to bqq'}, \; p_{H \to b\bar{b}}, \; p_{QCD}, \; \ldots$

# ongoing work: scaling this up

⊙ with the 128-dim "latent representation" of the jet calibrated,

✦ we observe that the original 10 classification scores close very well.

# ongoing work: scaling this up

⦿ with the 128-dim "latent representation" of the jet calibrated,

  ✦ we observe that the original 10 classification scores close very well.
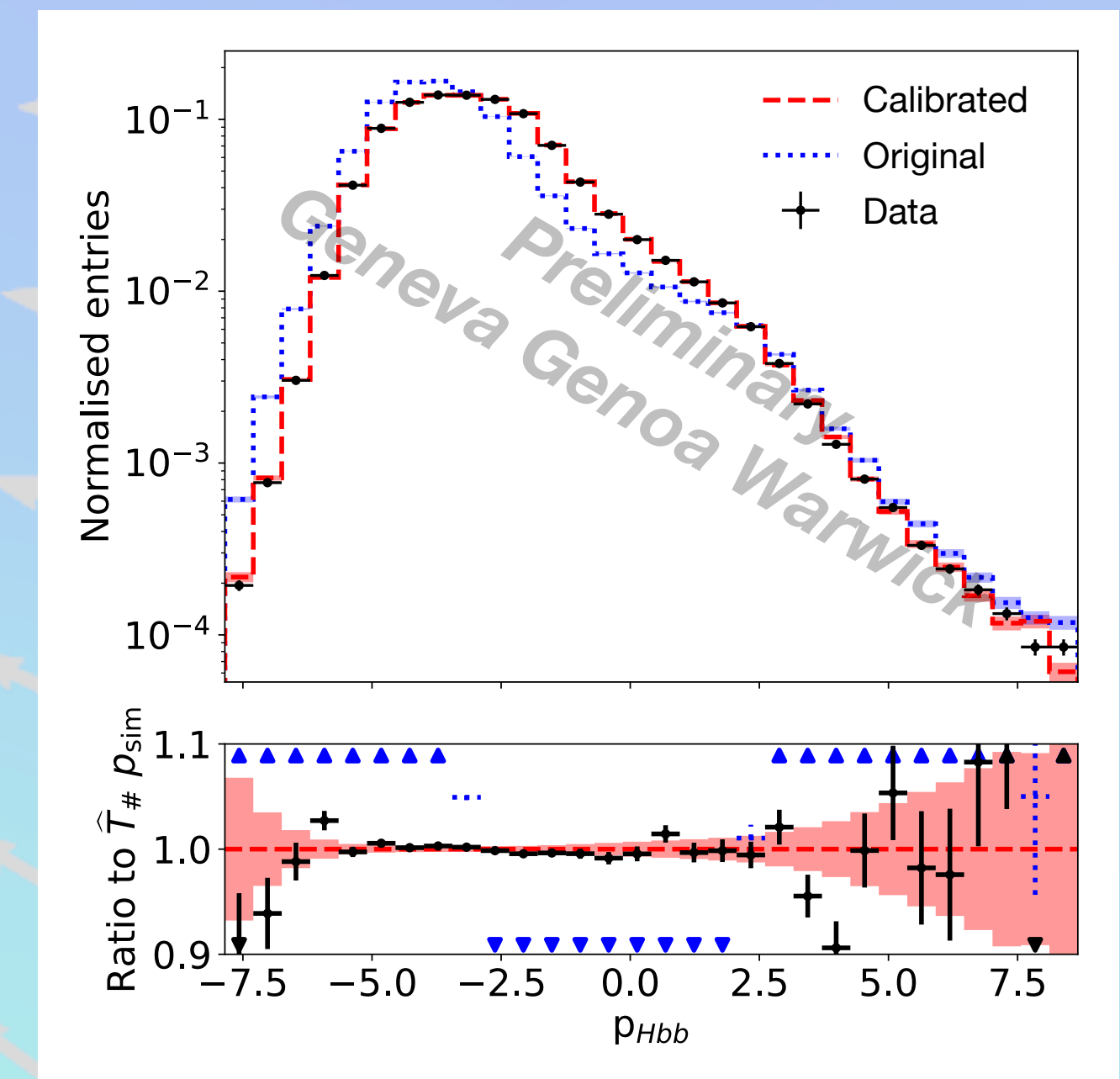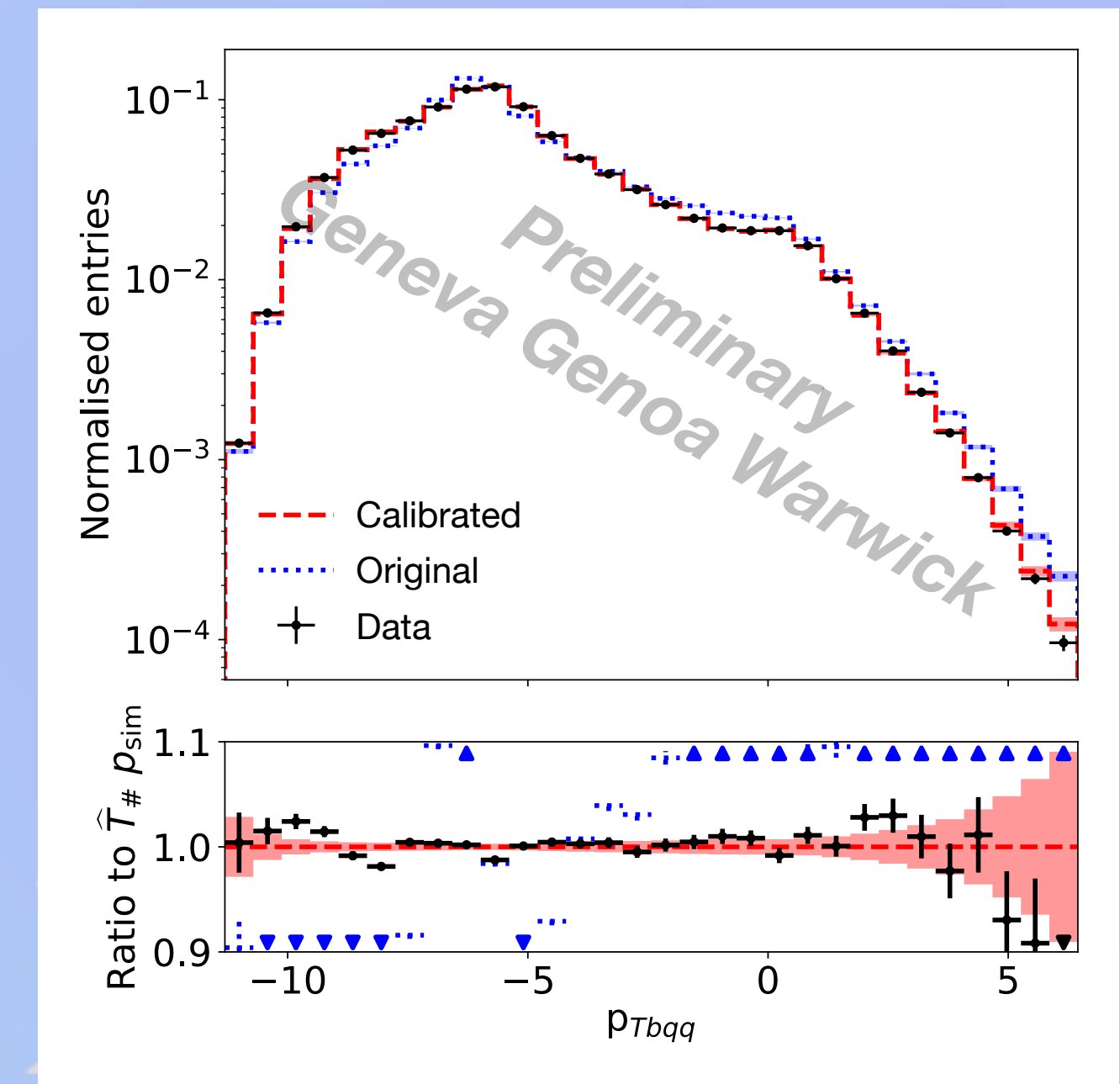
# ongoing work: scaling this up

- with the 128-dim "latent representation" of the jet calibrated,

  - ✦ we observe that the original 10 classification scores close very well.

- but this enables "arbitrary" use of the information contained in that representation:

particle inputs

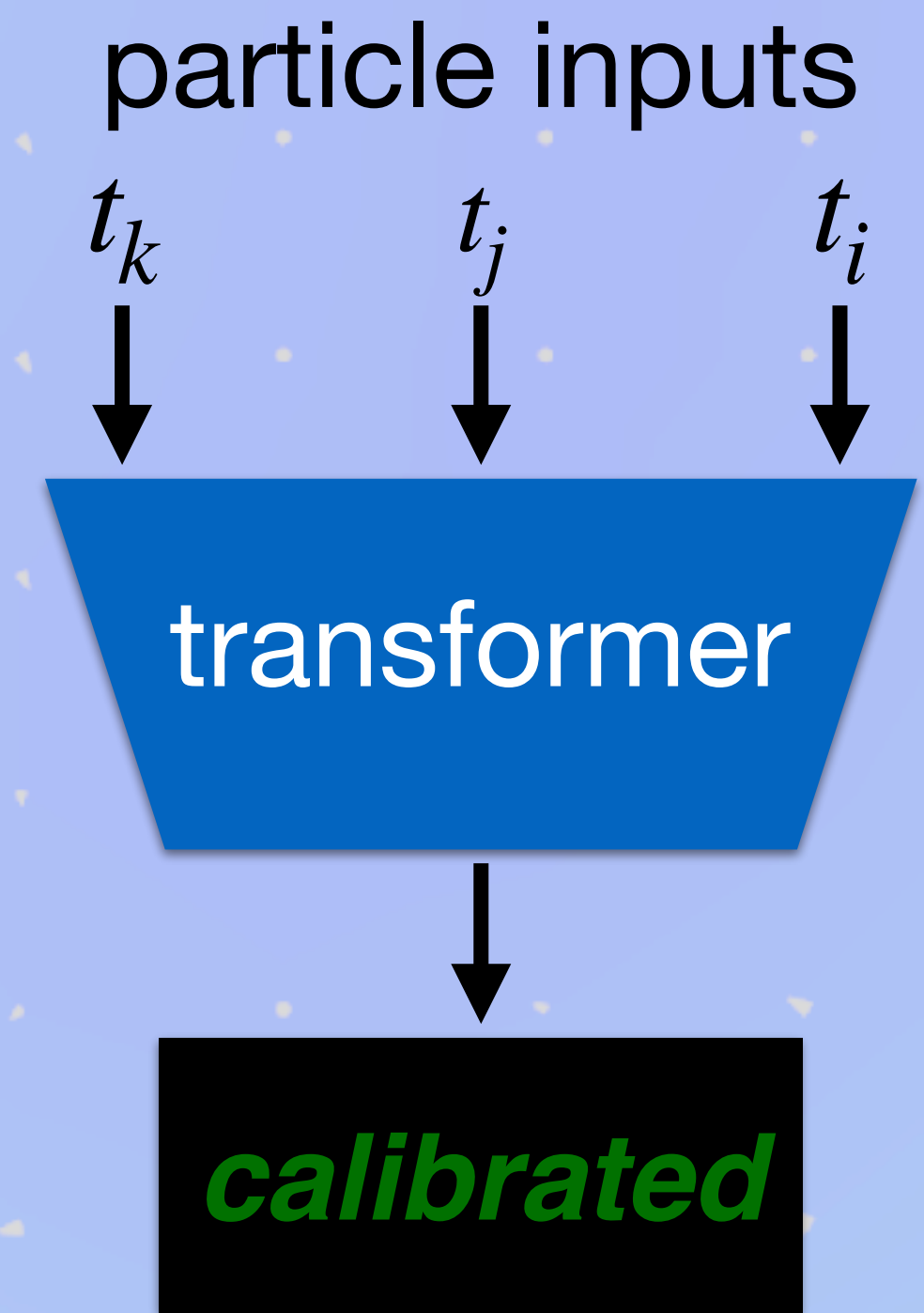$t_k$     $t_j$     $t_i$

transformer

*calibrated*

# ongoing work: scaling this up

◉ with the 128-dim "latent representation" of the jet calibrated,

   ✦ we observe that the original 10 classification scores close very well.

◉ but this enables "arbitrary" use of the information contained in that representation:

particle inputs

$t_k$     $t_j$     $t_i$

transformer

*calibrated*

your

new

MLP

"automatically" calibrated output

# ongoing work: scaling this up

⊙ with the 128-dim "latent representation" of the jet calibrated,

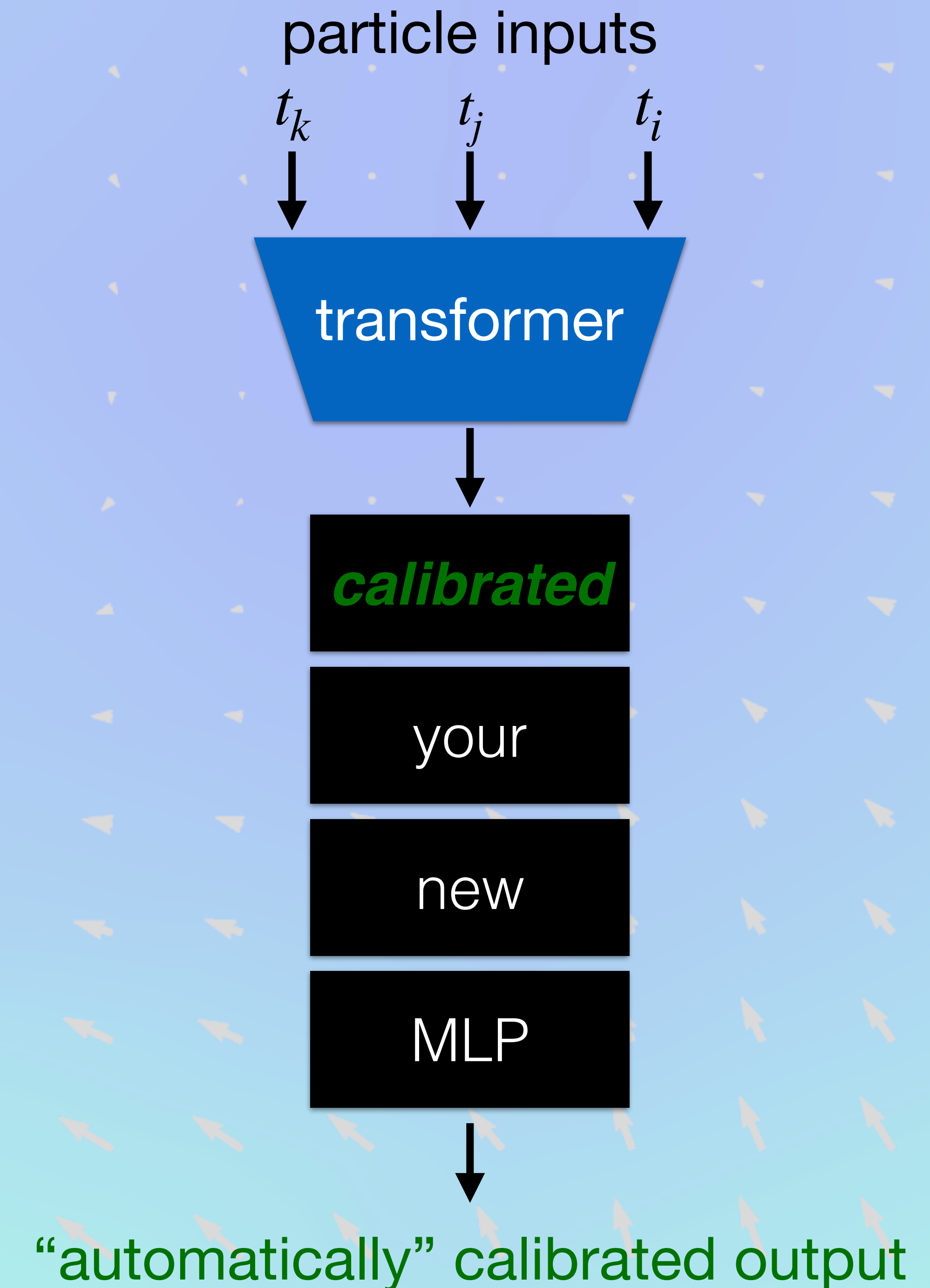✦ we observe that the original 10 classification scores close very well.

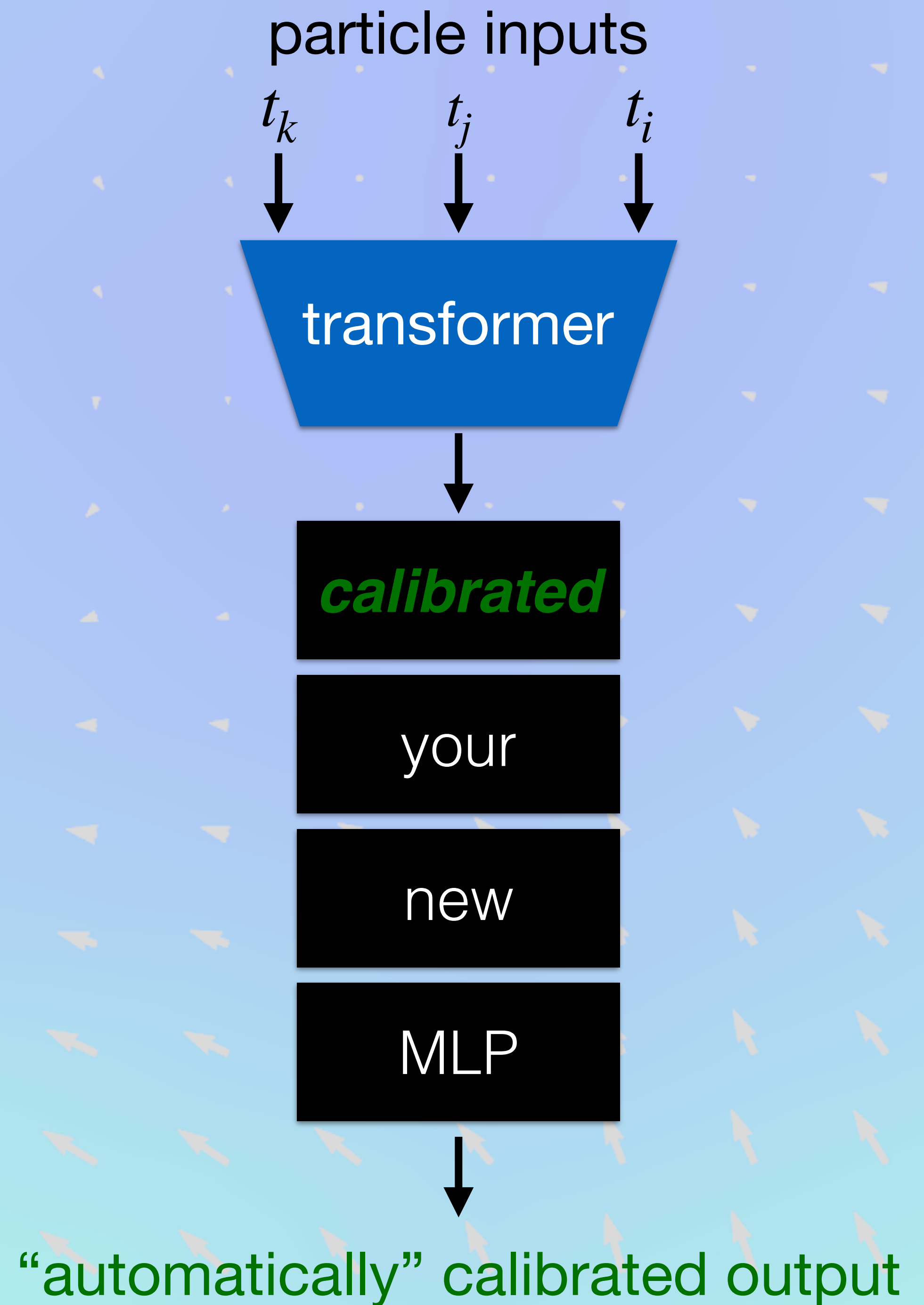⊙ but this enables "arbitrary" use of the information contained in that representation:

✦ allows calibration of "foundation" models for broad use.

✦ useful well beyond particle physics!

particle inputs

$t_k$      $t_j$      $t_i$

transformer

calibrated

your

new

MLP

"automatically" calibrated output

67

**nSBI promises to yield the best possible constraints on a parameter $\phi$ given an observable $x$.**

⊙ that's a **big** claim but also worth pursuing…

⊙ *all* the information of $x$ relevant to $\phi$ in the simulator $p(x \mid \phi)$ can be used for inference:

✦ allows more complex $x$ and $\phi$ than standard approaches.

✦ potential huge gains in cases where NPs are very correlated with PoIs.

✦ **requires very good (correct!) simulation.**

**nSBI promises to yield the best possible constraints on a parameter $\phi$ given an observable $x$.**

⦿ lots of ongoing work to make this a reality:

✦ to contend with "technical" problems like convergence and global NPs,

✦ to calibrate complicated observables $x$ to enable ***correct*** constraints.

⦿ large potential speed-up in inference time.

⦿ I'm hopeful that in $\mathcal{O}(5)$ years there will be "over the counter" solutions available — but it will take effort to get there.
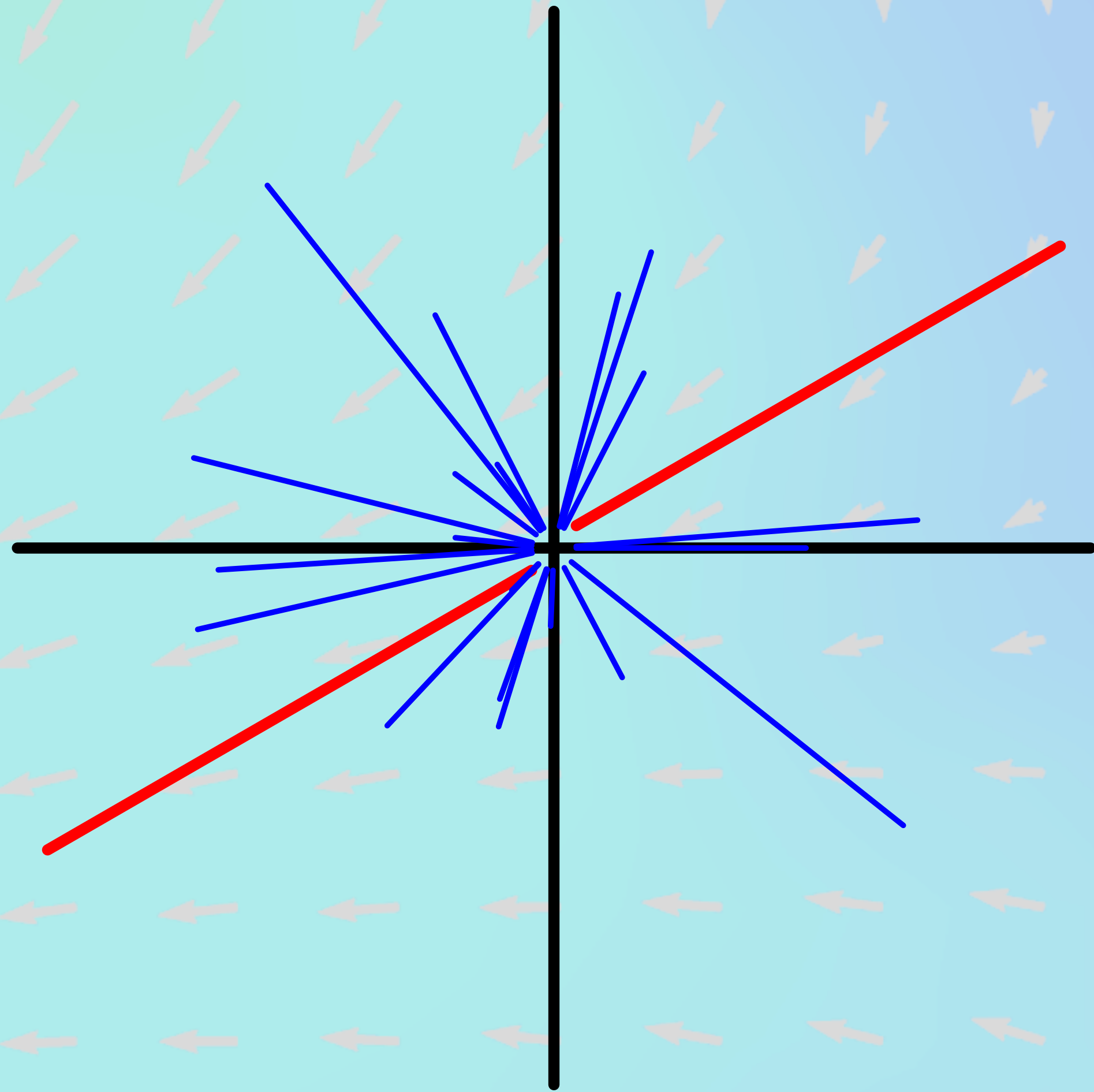
✦ It's a very interesting space with lots of work to do!

thank you!

# why do we need summaries / embeddings / "physics objects"?

😱 **tractability**: ~millions of detector channels to read out per LHC bunch crossing.

😱 **correctness**: difficult to construct a simulator that *adequately describes all details* of the data!

**indeed, this is perhaps the main "point" of constructing jets:**

◉ we cannot correctly predict the details of QCD with arbitrary accuracy;

◉ we *can* predict the "large-scale structure" of the fragmentation of partons.

tough to predict

summary

*"large-scale structure": calculable features*

**calibration region**

**analysis region**

**even so, having access *more relevant details* can enable better constraints.**

◉ we can often *measure* the density over some features better than we can predict it.

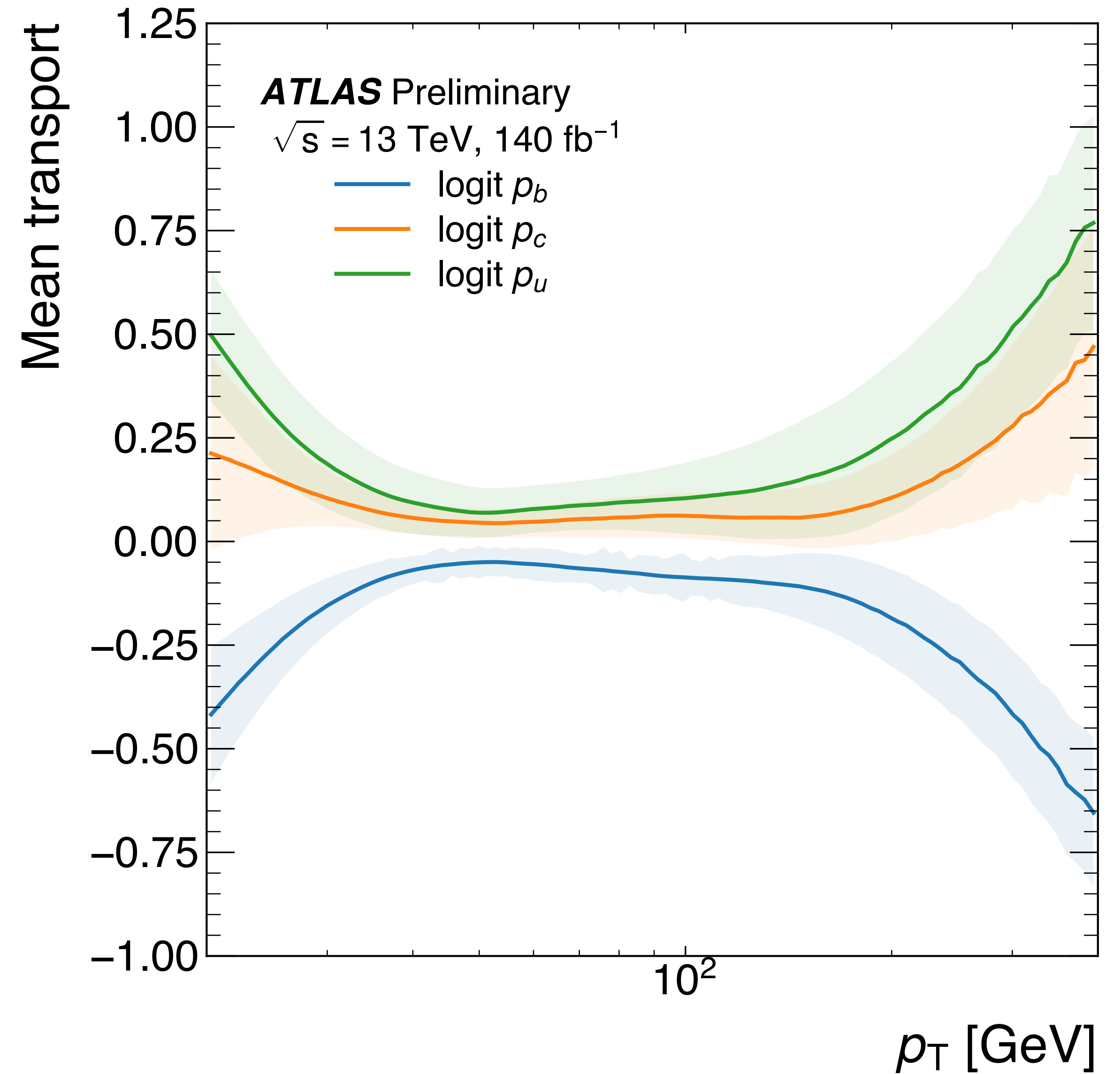◉ a *calibration region* may be used to measure said density and to *correct* the simulation.

notation:

- $q_i \equiv \text{logit}\, p_i$ : flav. class. scores
- $p_{\text{sim}}(\vec{q} \mid p_T) \equiv p_{\text{sim}}(\vec{q} \mid p_T)\, p_{\text{data}}(p_T)$
- $\hat{T}_\# \equiv p_T$-dependent OT map

◉ for $b$-jets: $\hat{T} p_b < p_b$, while the reverse is true for $p_c$ and $p_u$.

    ✦ the simulation overstates its classification power.

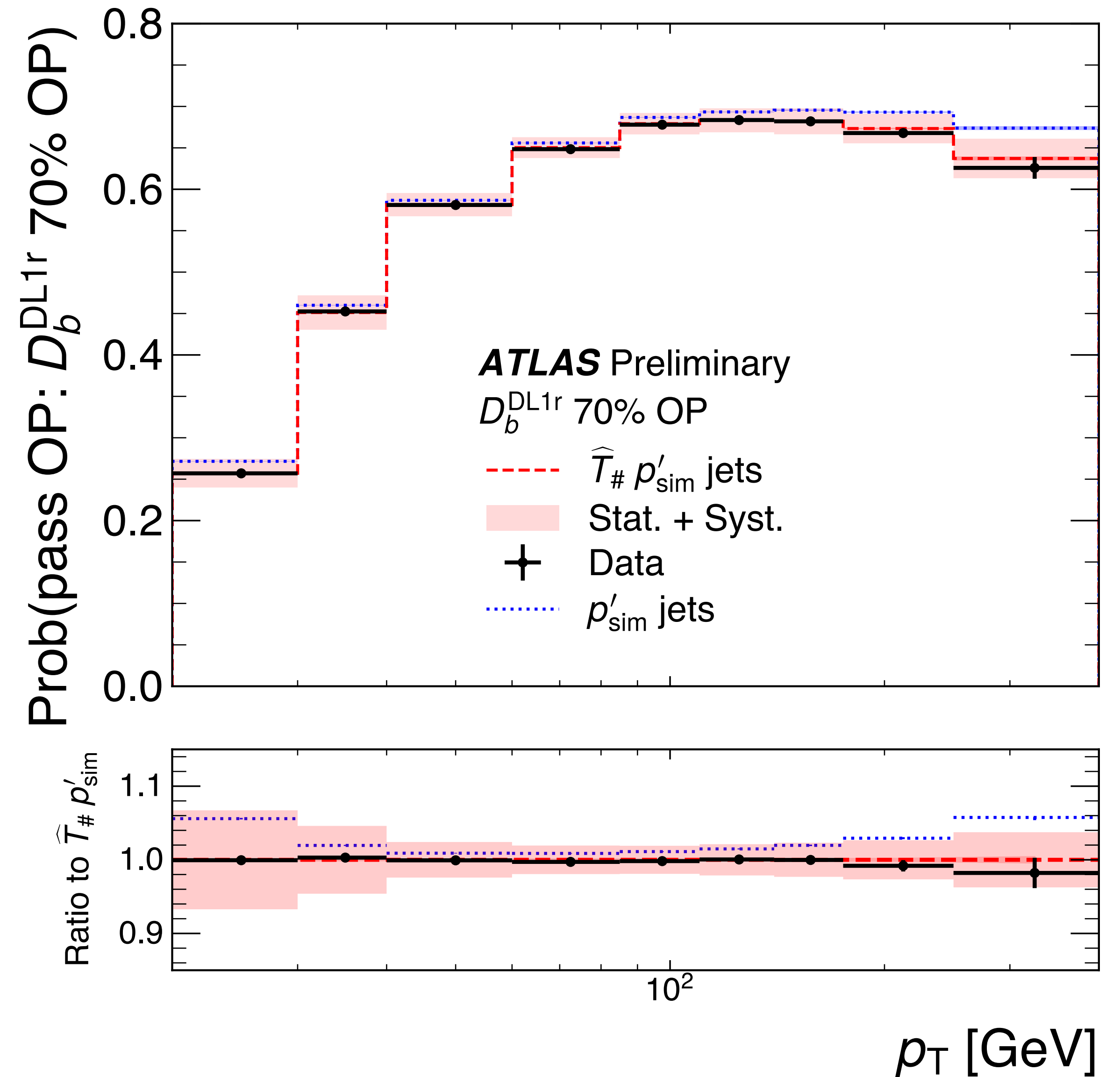◉ we have much more information about what aspects of the simulation are incorrect than before.



ATLAS Preliminary
$\sqrt{s}$ = 13 TeV, 140 fb$^{-1}$
— logit $p_b$
— logit $p_c$
— logit $p_u$

Mean transport

$p_T$ [GeV]

notation:

- $q_i \equiv \mathrm{logit}\, p_i$ : flav. class. scores
- $p_{\mathrm{sim}}(\vec{q}\,|\,p_T) \equiv p_{\mathrm{sim}}(\vec{q}\,|\,p_T)\, p_{\mathrm{data}}(p_T)$
- $\hat{T}_{\#} \equiv p_T$-dependent OT map

**conventional operating points are "automatically" corrected.**

⊙ excellent closure observed for the "standard" $b$-tagging discriminant points

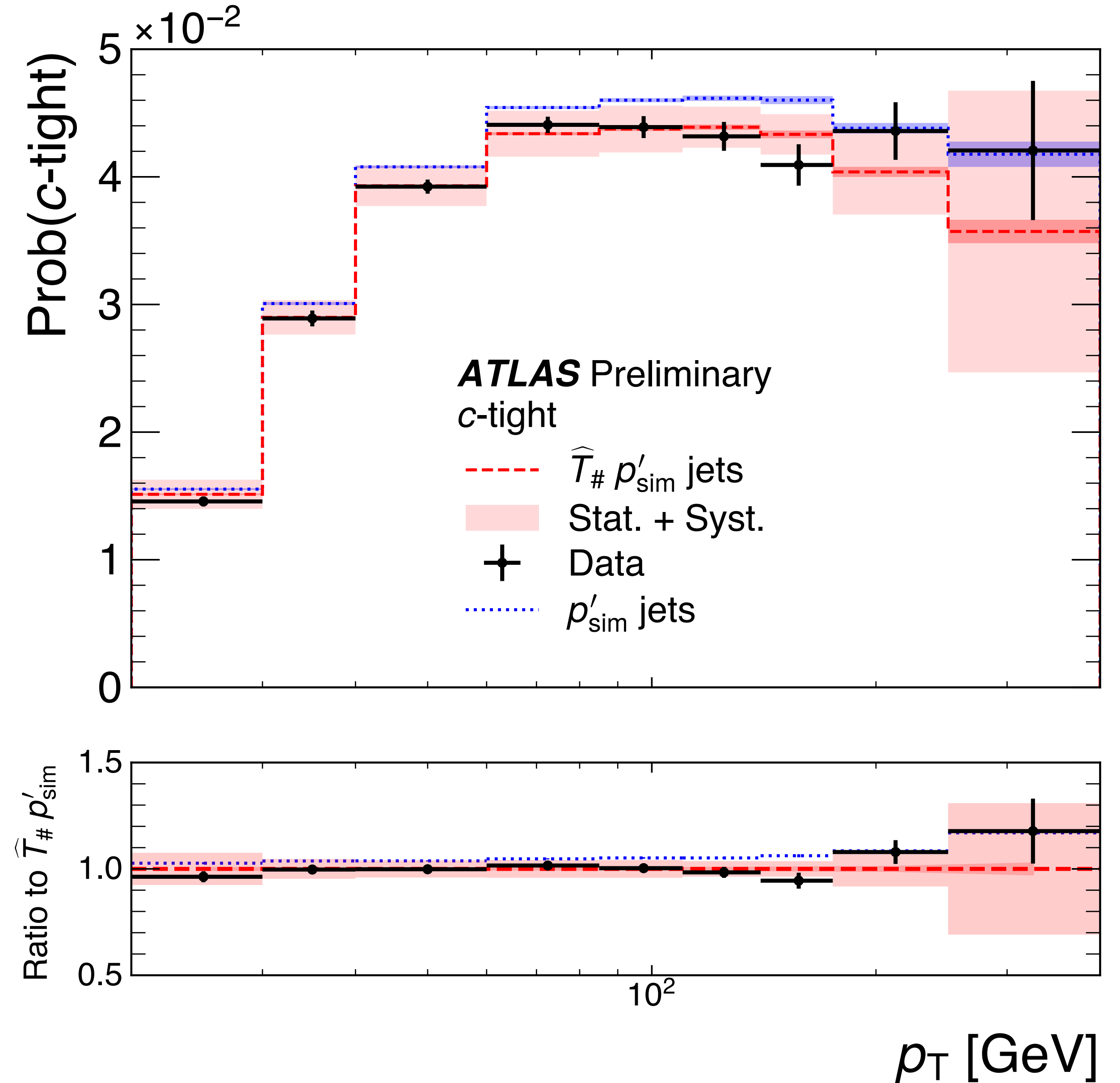✦ used to observe Higgs couplings to $W, b, t$.



**ATLAS** Preliminary
$D_b^{\mathrm{DL1r}}$ 70% OP

- - - $\hat{T}_{\#}\, p'_{\mathrm{sim}}$ jets
  Stat. + Syst.
  ＋ Data
  ⋯⋯ $p'_{\mathrm{sim}}$ jets

Prob(pass OP: $D_b^{\mathrm{DL1r}}$ 70% OP)

Ratio to $\hat{T}_{\#}\, p'_{\mathrm{sim}}$

$p_{\mathrm{T}}$ [GeV]

notation:

- $q_i \equiv \operatorname{logit} p_i$ : flav. class. scores
- $p_{\mathrm{sim}}(\vec{q}\,|\,p_T) \equiv p_{\mathrm{sim}}(\vec{q}\,|\,p_T)\,p_{\mathrm{data}}(p_T)$
- $\hat{T}_{\#} \equiv p_T$-dependent OT map

**more general uses of flavor-tagging become possible.**

◉ charm-tagging discriminator used to constrain the $H \leftrightarrow c$ couplings shows good agreement.

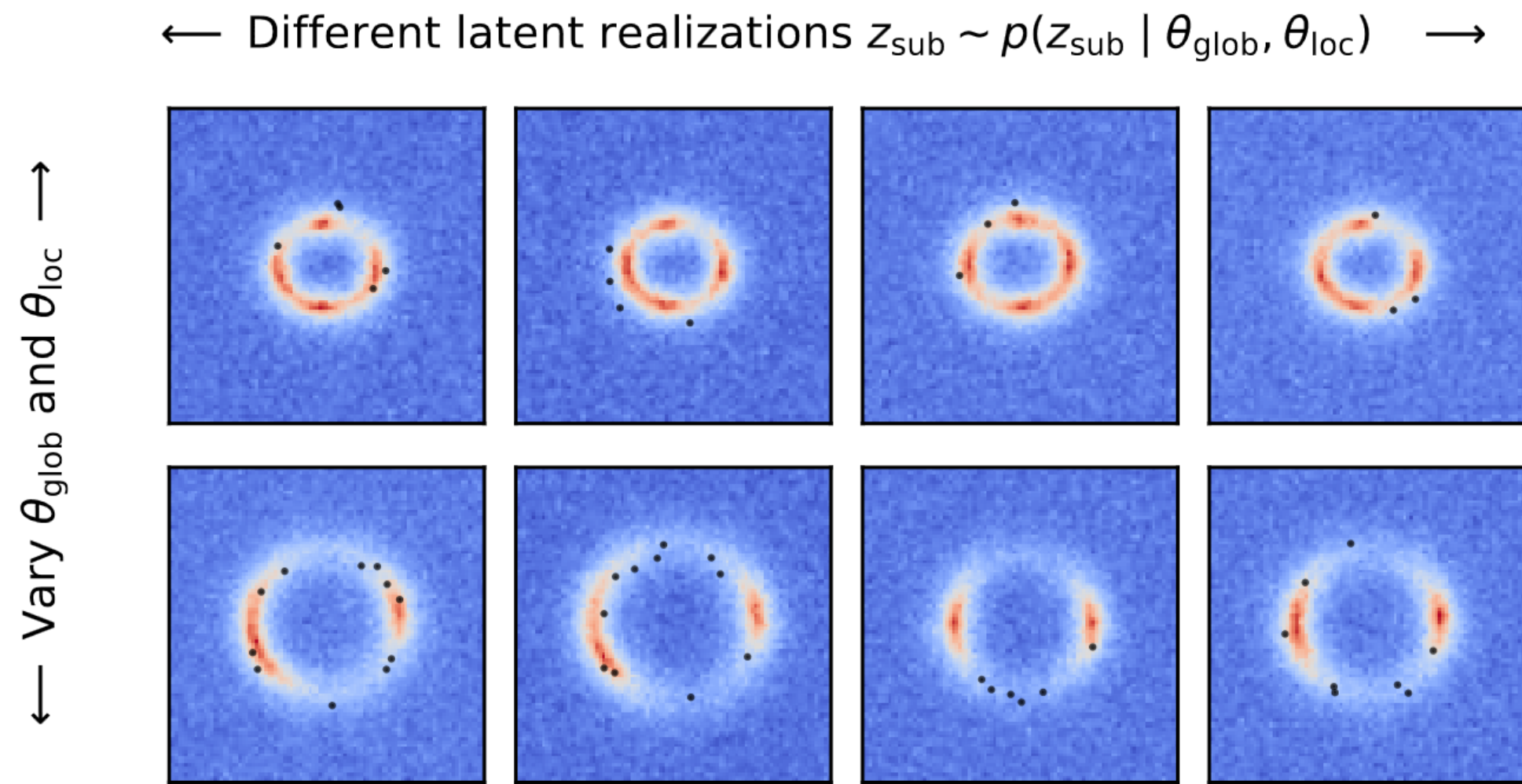◉ closure depends on the full 3D density $\hat{T}_{\#}\, p_{\mathrm{sim}}(\vec{q}\,|\,p_T)$ agreeing with data.

← Different latent realizations $z_{sub} \sim p(z_{sub} \mid \theta_{glob}, \theta_{loc})$ →
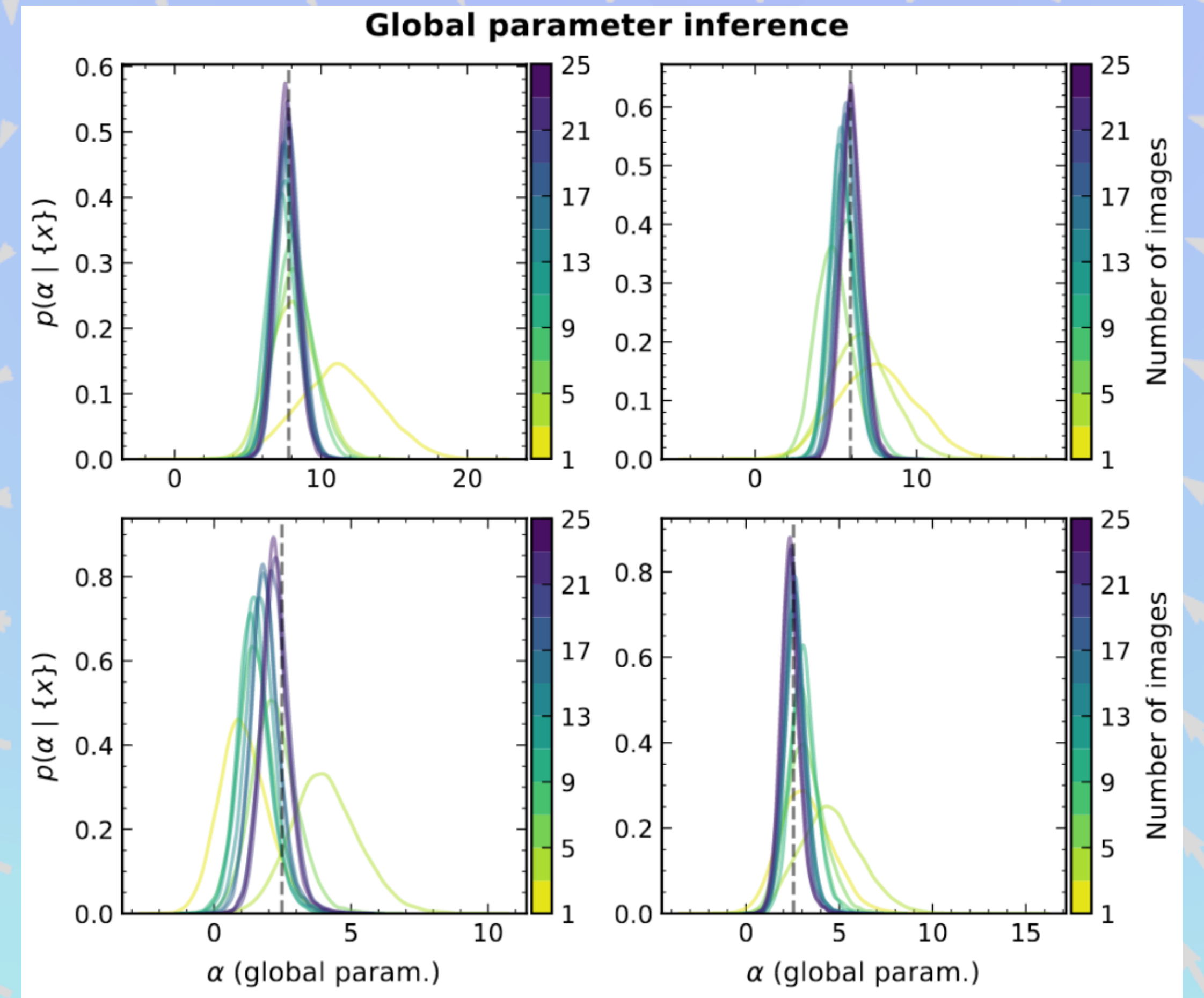
Vary $\theta_{glob}$ and $\theta_{loc}$

Figure 5: Illustrative samples from the lensing model. The rows show two different choices of local (per-event) parameters, while the different columns show variations on the global (set-wide) parameters for the fixed choice of local parameters. Sample-to-sample variation induced by the global parameters, which control the abundance of a subhalo population in the lens, can be seen. The scatter points shows the location of individual subhalos in each image.
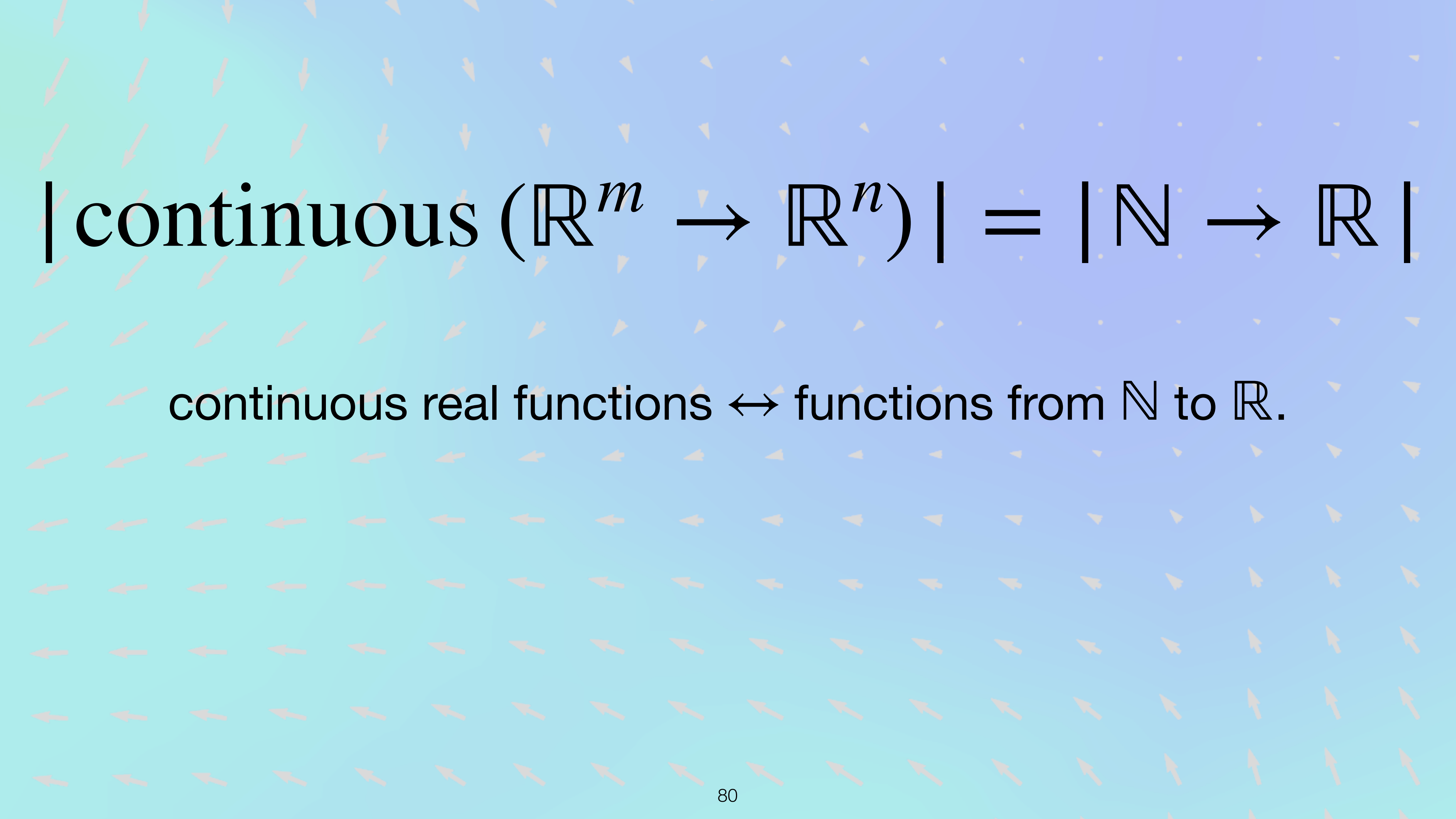


Global parameter inference

credit: Siddarth Mishra-Sharma

before we can talk about neural inference,

we have to believe machine-learning is doing what it claims!

why does it work?
(from the perspective of a physicist…)

$$|\text{continuous}\,(\mathbb{R}^m \to \mathbb{R}^n)| = |\mathbb{N} \to \mathbb{R}|$$

continuous real functions $\leftrightarrow$ functions from $\mathbb{N}$ to $\mathbb{R}$.

$$|\text{continuous}\,(\mathbb{R}^m \to \mathbb{R}^n)| = |\mathbb{N} \to \mathbb{R}|$$

continuous real functions $\leftrightarrow$ functions from $\mathbb{N}$ to $\mathbb{R}$.

i.e. any continuous function can be *parameterized* by a countable number of real numbers.
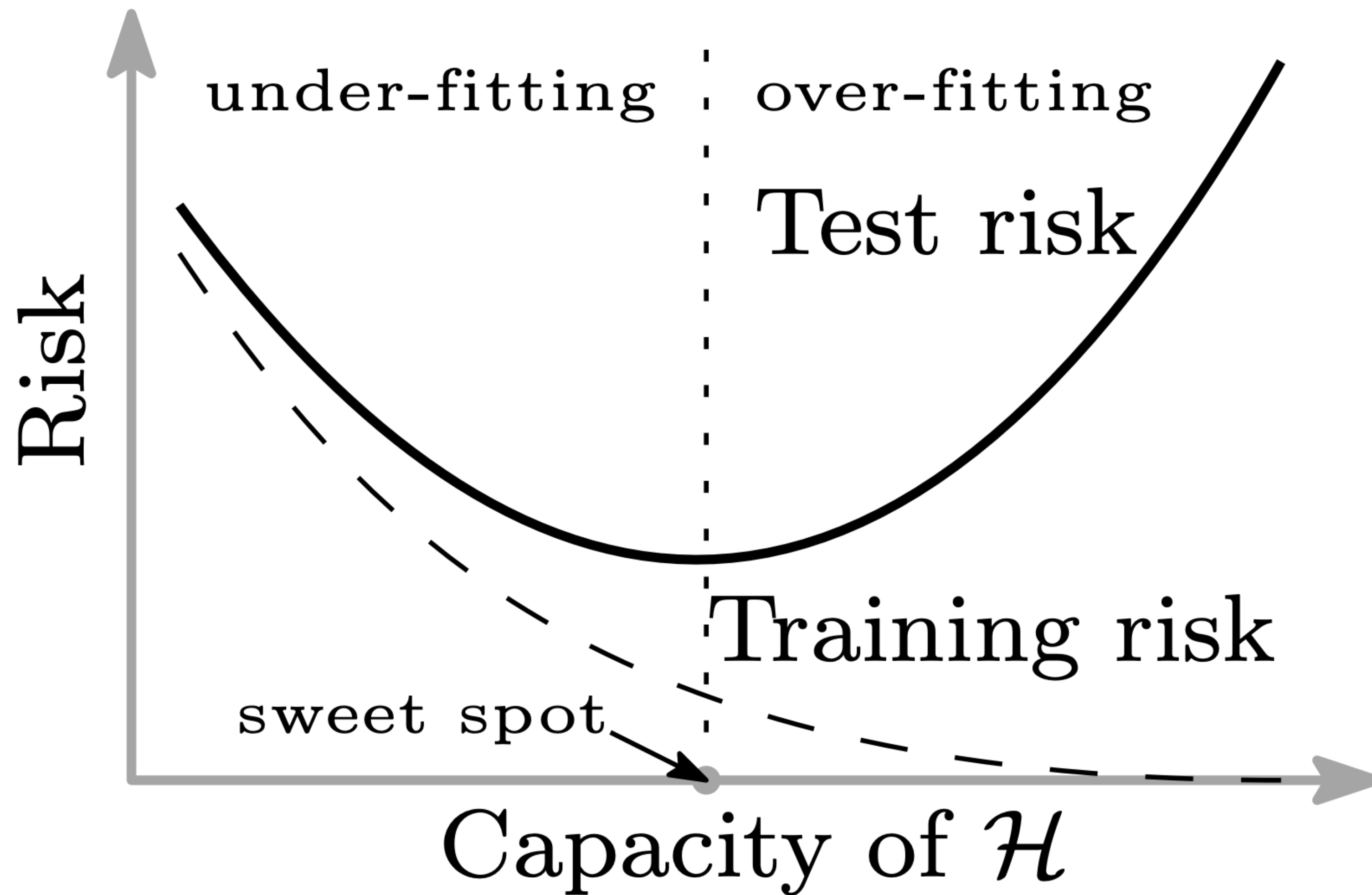
$$|\text{continuous}\,(\mathbb{R}^m \to \mathbb{R}^n)| = |\mathbb{N} \to \mathbb{R}|$$

continuous real functions $\leftrightarrow$ functions from $\mathbb{N}$ to $\mathbb{R}$.

i.e. any continuous function can be *parameterized*
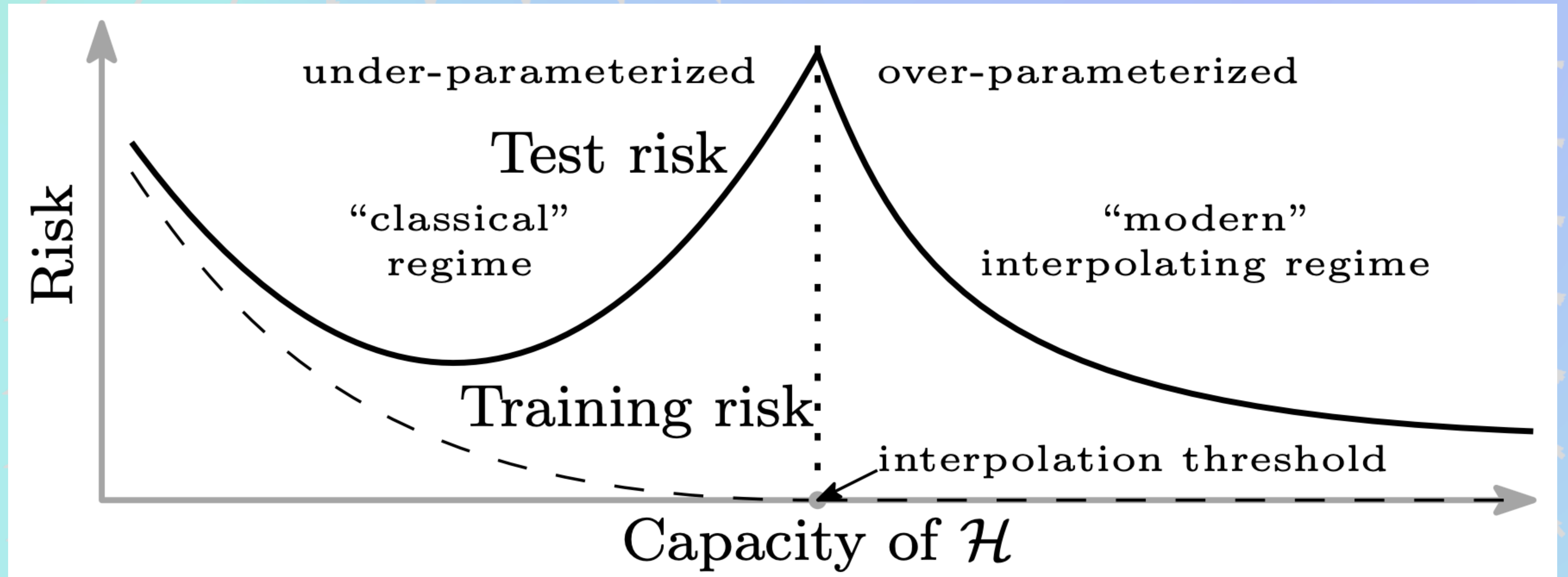by a countable number of real numbers.

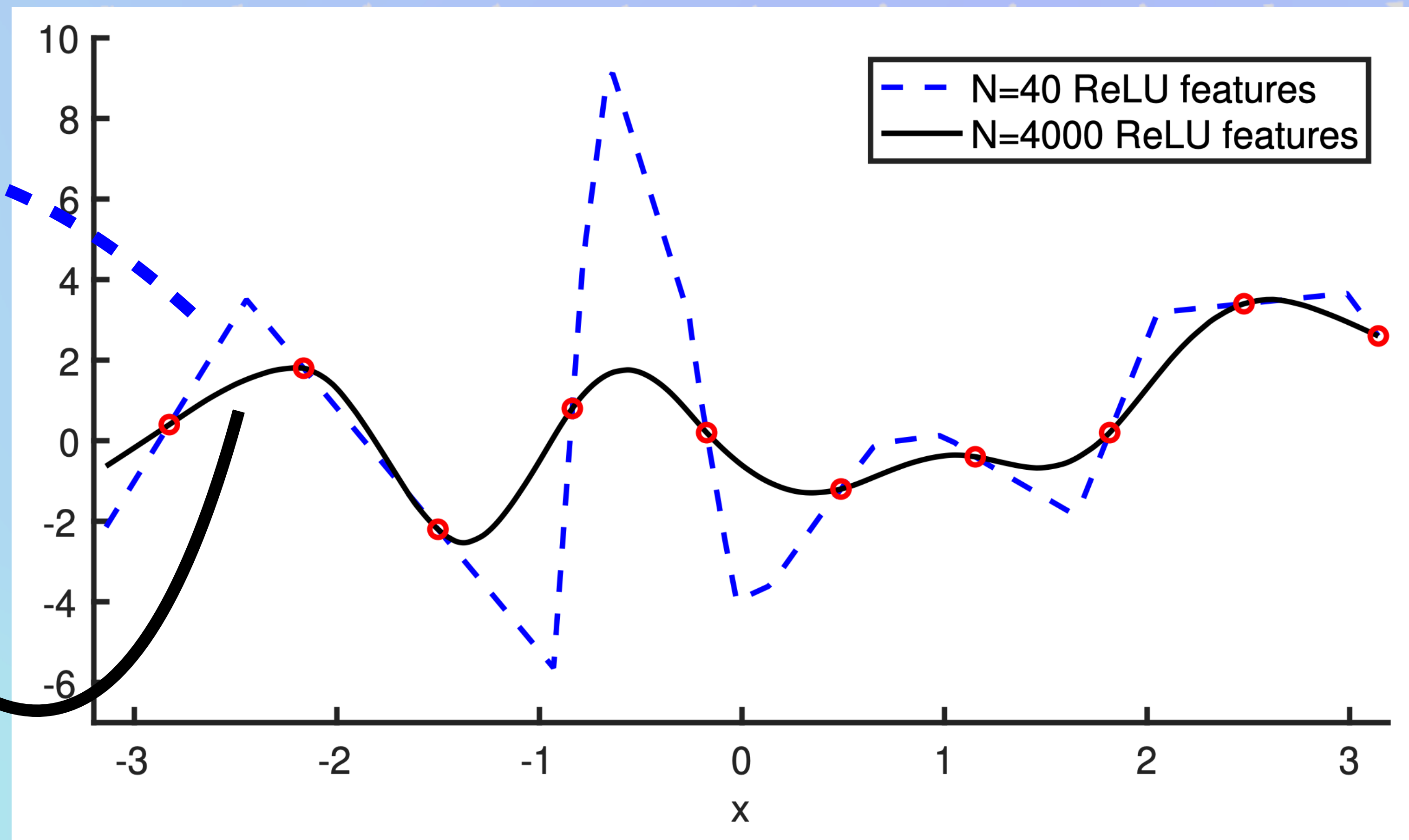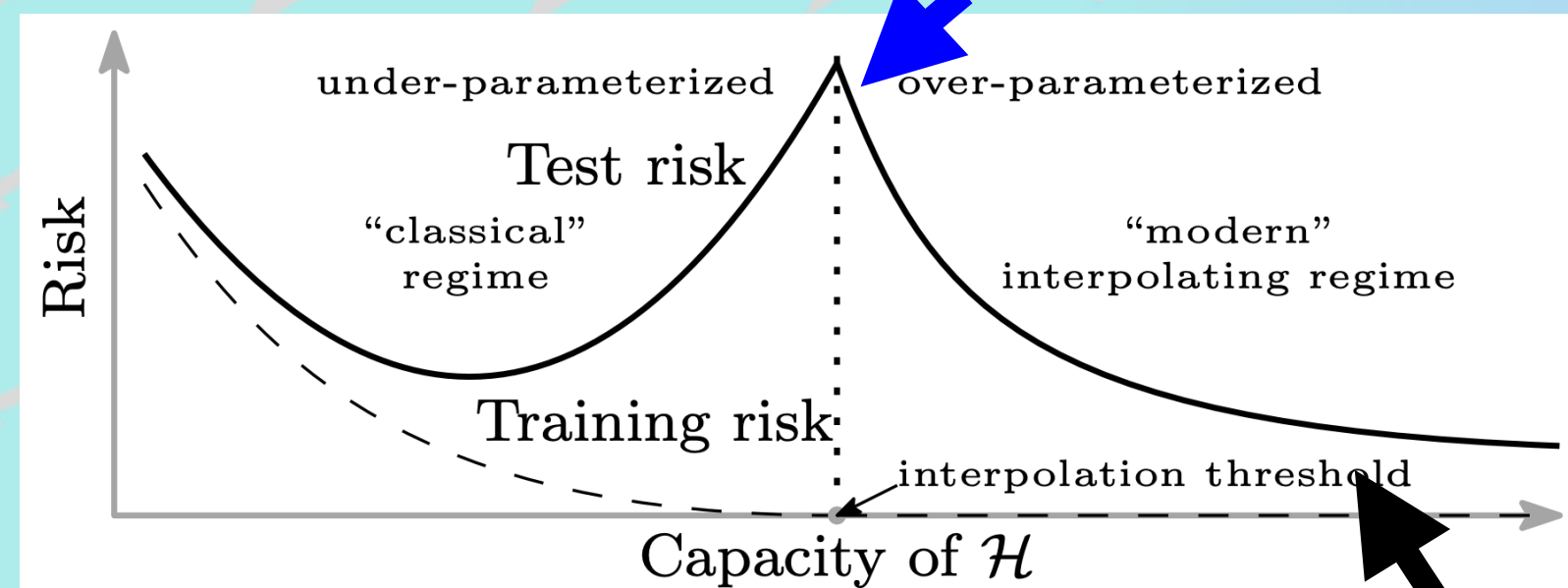(underpins e.g. Fourier series, eigenvector decomposition, etc)

"classical" view of fitting data

Risk ~ error of fit on unseen data

growing (empirical) evidence:
"classical" under- and over-fitting ideals do not apply.

Test risk / Training risk diagram: under-parameterized "classical" regime, over-parameterized "modern" interpolating regime, interpolation threshold, Capacity of $\mathcal{H}$. Plot of N=40 ReLU features (dashed blue) and N=4000 ReLU features (solid black).
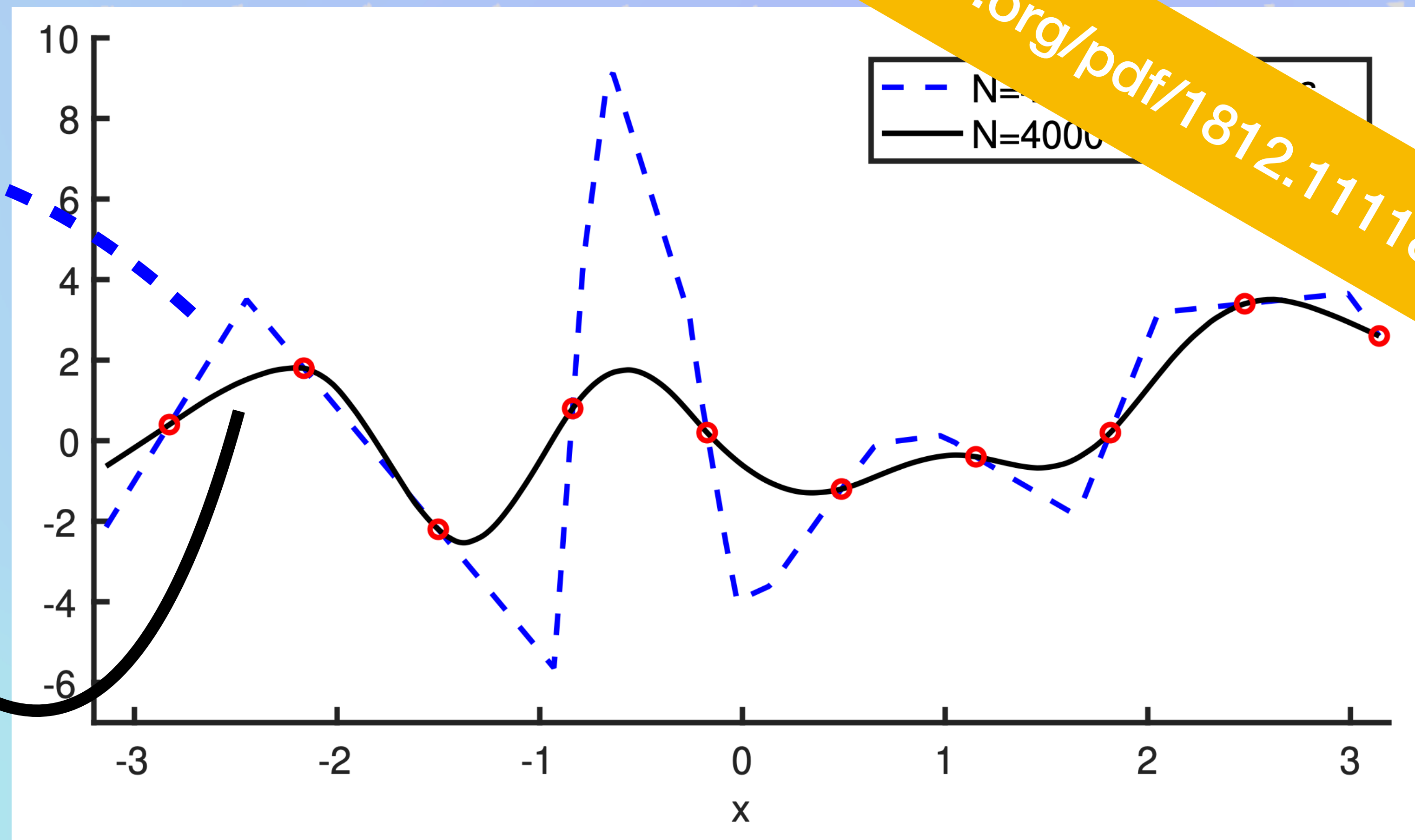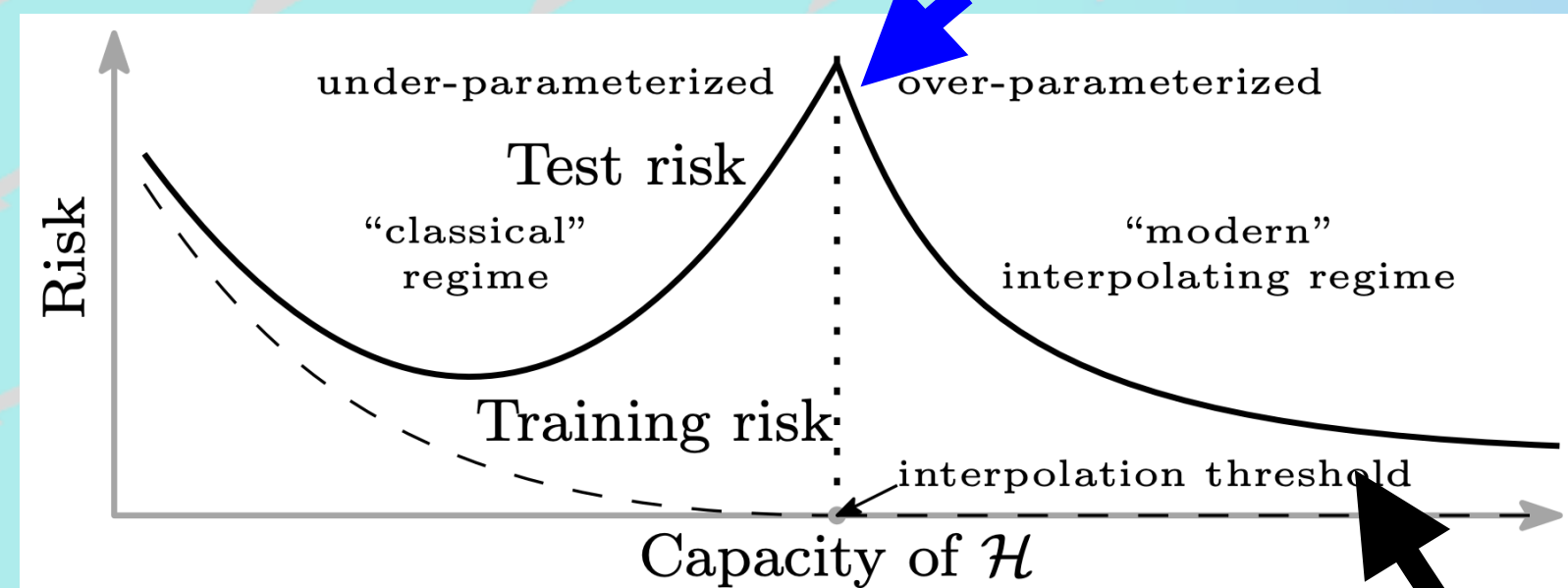
growing (empirical) evidence:
"classical" under- and over-fitting ideals do not apply.

growing (empirical) evidence:
"classical" under- and over-fitting ideals do not apply.

to summarize:

1) NNs approximate any real function to arbitrary precision
with a finite number of parameters (and enough training data).

2) over-parameterization yields high-quality (very predictive) fits.

3) over-parameterization *increases* the odds the fit converges to a
"good" local minimum.

mounting evidence that NNs really are learning what they claim.

to summarize:

1) NNs approximate any real function to arbitrary precision
with a finite number of parameters (and enough training data).

2) over-parameterization yields high-quality (very predictive) fits.

$\rightarrow$ growing confidence that NNs are learning what they claim to.