

Reps4LHC

Tilman Plehn

Symmetries

Uncertainties

Explainability

Anomalies

Representation Learning for the LHC

Tilman Plehn

Universität Heidelberg

Annual CRC Meeting, September 2025



ML as representation learning

Similar to fit

- approximate $f_\theta(x) \approx f(x)$
 - x phase space
 - f_θ numerical function
- θ data representation

Phase space probabilities

- regression $x \rightarrow A_\theta(x)$
- classification $x \rightarrow p_\theta(x)$ [likelihood ratio]
- generation $r \sim \mathcal{N} \rightarrow x \sim p_\theta(x)$
- conditional generation $r \sim \mathcal{N} \rightarrow x \sim p_\theta(x|y)$

Requirements on θ

- accuracy
 - precision
 - structure [physics?]
- Physics knowledge or knowledge-free upscaling?



Lorentz-equivariance

Encode known symmetries

- permutation co-variance → graph or transformer
 - Lorentz co-variance → $\Lambda(f_\theta(x)) = f_\theta(\Lambda(x))$ [4-vectors vs Mandelstams]
- 1- L-GATr geometric algebra representation
 - 2- LLoCa local reference frame per particle



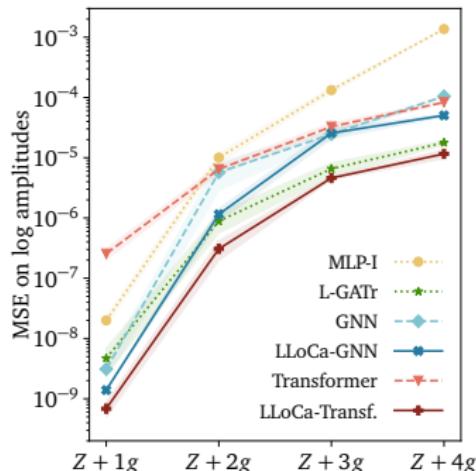
Lorentz-equivariance

Encode known symmetries

- permutation co-variance → graph or transformer
 - Lorentz co-variance → $\Lambda(f_\theta(x)) = f_\theta(\Lambda(x))$ [4-vectors vs Mandelstams]
- 1- L-GATr geometric algebra representation
 - 2- LLoCa local reference frame per particle

Performance is all you need

- LO transition amplitudes $q\bar{q} \rightarrow Z + 1\dots4 g$
- improved scaling with multiplicity



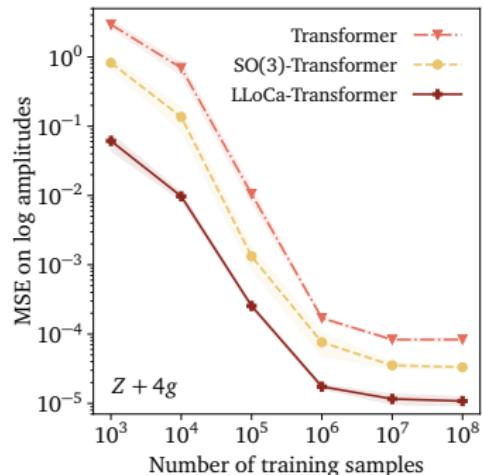
Lorentz-equivariance

Encode known symmetries

- permutation co-variance → graph or transformer
 - Lorentz co-variance → $\Lambda(f_\theta(x)) = f_\theta(\Lambda(x))$ [4-vectors vs Mandelstams]
- 1- L-GATr geometric algebra representation
 - 2- LLoCa local reference frame per particle

Performance is all you need

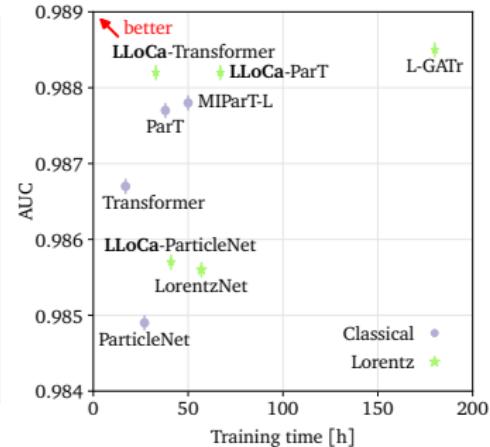
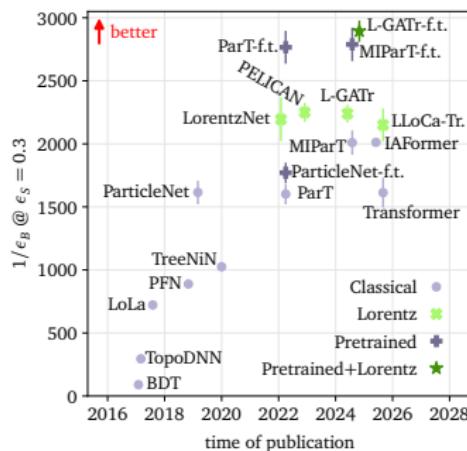
- LO transition amplitudes $q\bar{q} \rightarrow Z + 1\dots4 g$
 - improved scaling with multiplicity
 - subset of symmetries
- Advantage across implementations



Equivariant jet tagging

Tagging benchmarks

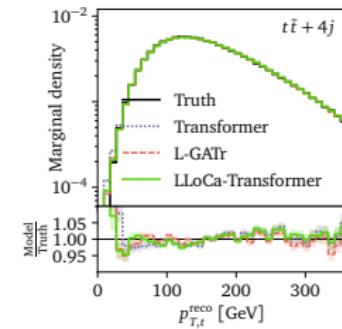
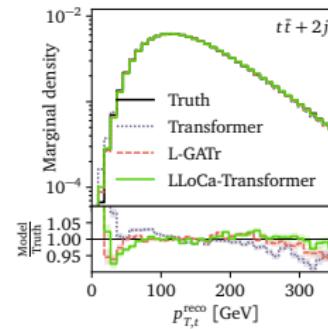
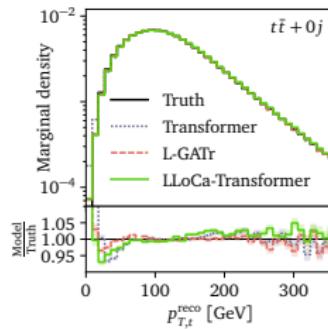
- flavor tagging and top tagging
→ more than 10x improvement over best BDT
 - problem: jet axis breaking Lorentz-symmetry
 - 1- give jet axis explicitly as additional 4-vector
 - 2- reduce encoded symmetry [LLoCa]
 - multi-class tagging the same
- Resilience, uncertainties, experiment next



Equivariant event generation

MadGraph7: ML-event generation

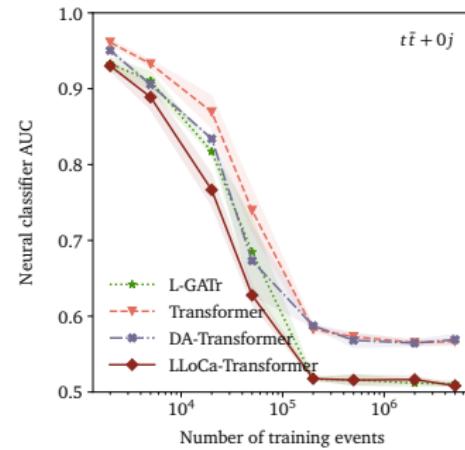
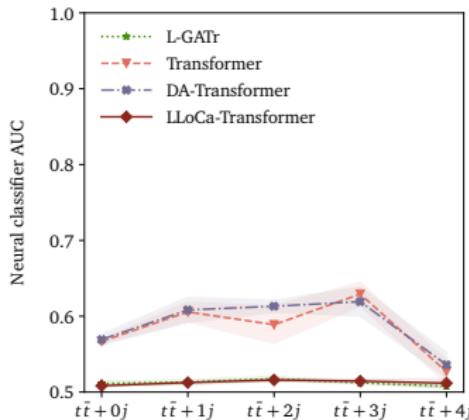
- conditional flow matching with transformer-velocity
- end-to-end $pp \rightarrow t_{\text{had}}\bar{t}_{\text{had}} + 0\dots4j$ [22M training events]
- per-cent kinematics



Equivariant event generation

MadGraph7: ML-event generation

- conditional flow matching with transformer-velocity
 - end-to-end $p\bar{p} \rightarrow t_{\text{had}}\bar{t}_{\text{had}} + 0\dots 4j$ [22M training events]
 - per-cent kinematics
 - LHC: training-generation classifier
- Accurate phase space a solved problem...



Representing mean and uncertainty

Remember a fit

- learn scalar field $f_\theta(x) \approx f(x)$
- statistics: maximize parameter probability given (f_j, σ_j)

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

→ maximize likelihood instead

$$\begin{aligned} p(x|\theta) &= \prod_j \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{|f_j - f_\theta(x_j)|^2}{2\sigma_j^2}\right) \\ \Rightarrow \quad \mathcal{L} \equiv -\log p(x|\theta) &= \sum_j \frac{|f_j - f_\theta(x_j)|^2}{2\sigma_j^2} + \text{const}(\theta) \end{aligned}$$



Representing mean and uncertainty

Remember a fit

- learn scalar field $f_\theta(x) \approx f(x)$
- statistics: maximize parameter probability given (f_j, σ_j)

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

→ maximize likelihood instead

$$\begin{aligned} p(x|\theta) &= \prod_j \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{|f_j - f_\theta(x_j)|^2}{2\sigma_j^2}\right) \\ \Rightarrow \quad \mathcal{L} &\equiv -\log p(x|\theta) = \sum_j \frac{|f_j - f_\theta(x_j)|^2}{2\sigma_j^2} + \text{const}(\theta) \end{aligned}$$

Learned local uncertainty

- Gaussian log-likelihood with normalization

$$\mathcal{L}_{\text{heteroscedastic}} = \frac{|f(x) - f_\theta(x)|^2}{2\sigma_\theta(x)^2} + \log \sigma_\theta(x) + \dots$$

- if needed replace $\sigma_\theta(x)$ by mixture model
- learn $f_\theta(x)$ and $\sigma_\theta(x)$ together



1- Bayesian networks

Learned function statistically

- amplitude over phase phase

$$\langle A \rangle = \int dA \ A \ p(A)$$

- internal representation θ of training data T [think Gaussian with mean and width]

$$p(A) = \int d\theta \ p(A|\theta) \ p(\theta|T)$$

→ θ -distribution defining Bayesian NN

Variational approximation

- definition of training

$$p(A) = \int d\theta \ p(A|\theta) \ p(\theta|T) \approx \int d\theta \ p(A|\theta) \ q(\theta)$$

- similarity through minimal KL-divergence [Bayes' theorem to remove unknown posterior]

$$\begin{aligned} D_{KL}[q(\theta), p(\theta|T)] &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)}{p(\theta|T)} \\ &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)p(T)}{p(T|\theta)p(\theta)} \\ &\approx D_{KL}[q(\theta), p(\theta)] - \int d\theta \ q(\theta) \ \log p(T|\theta) \equiv \mathcal{L} \end{aligned}$$

→ Two-term loss: likelihood + prior



1- Bayesian networks

Learned function statistically

- amplitude over phase phase

$$\langle A \rangle = \int dA \, A \, p(A)$$

- internal representation θ of training data T [think Gaussian with mean and width]

$$p(A) = \int d\theta \, p(A|\theta) \, p(\theta|T)$$

→ θ -distribution defining Bayesian NN

Two uncertainties

- statistical — vanishing for perfect training: $q(\theta) \rightarrow \delta(\theta - \theta_0)$

$$\sigma_{\text{stat}}^2 = \int d\theta \, q(\theta) \left[\bar{A}(\theta) - \langle A \rangle \right]^2$$

- systematic — vanishing for perfect data: $p(A|\theta) \rightarrow \delta(A - A_0)$

$$\sigma_{\text{syst}}^2 = \int d\theta \, q(\theta) \left[\bar{A^2}(\theta) - \bar{A}(\theta)^2 \right]$$

→ Systematics dominant for LHC



2- Repulsive ensembles

Posterior from network ensemble

- OED vs continuity equation

$$\frac{d\theta}{dt} = v(\theta, t) \quad \Leftrightarrow \quad \frac{\partial \rho(\theta, t)}{\partial t} = -\nabla_\theta [v(\theta, t)\rho(\theta, t)]$$

- Fokker-Planck equation with stationary $\rho(\theta, t) = \pi(\theta)$

$$\frac{d\theta}{dt} = -\nabla_\theta \log \frac{\rho(\theta, t)}{\pi(\theta)}$$

- ODE describing training progress

$$\begin{aligned} \theta^{t+1} - \theta^t &\propto -\nabla_{\theta^t} \left[\log \rho(\theta^t) - \log \pi(\theta^t) \right] \\ &= -\nabla_{\theta^t} \left[\log \sum_j k(\theta^t, \theta_j^t) - \log p(\theta | x_{\text{train}}^t) \right] \equiv -\nabla_{\theta^t} \mathcal{L}_{\text{RE}} \end{aligned}$$

→ Joint ensemble training

Repulsive ensembles

- train network ensemble
- apply repulsive force kernel in function space

→ Alternative for statistical uncertainty



3- Evidential regression

Uncertainties from latent distribution, without sampling [Bahl, Elmer, TP, Winterhalder]

- evidential distribution

$$p(A) = \int d\lambda p(A|\lambda) p(\lambda|T) \approx \int d\lambda p(A|\lambda) p(\lambda|\theta_0)$$

assuming $p(A|\lambda) = \mathcal{N}(A|\bar{A}, \sigma^2)$ with $\lambda \equiv (\bar{A}, \sigma^2)$

- choose $p(\lambda|\theta_0)$ as conjugate prior

$$p(\lambda|\theta_0) = \frac{\beta^\alpha \sqrt{\nu}}{\Gamma(\alpha)\sqrt{2\pi\sigma^2}} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left(-\frac{2\beta + \nu(\gamma - \bar{A})^2}{2\sigma^2}\right)$$

with $\{\gamma, \nu, \alpha, \beta\} (x, \theta_0)$

- analytic likelihood: Student-t

$$p(A) = \text{St}\left(A \middle| \gamma, \frac{\beta(1+\nu)}{\nu\alpha}, 2\alpha\right)$$

$$A_{\text{NN}} = \int d\lambda \bar{A} p(\lambda|\theta_0) = \gamma$$

$$\sigma_{\text{syst}}^2 = \int d\lambda \sigma^2 p(\lambda|\theta_0) = \frac{\beta}{\alpha-1}$$

$$\sigma_{\text{stat}}^2 = \int d\lambda [\bar{A} - A_{\text{NN}}]^2 p(\lambda|\theta_0) = \frac{\beta}{\nu(\alpha-1)}$$

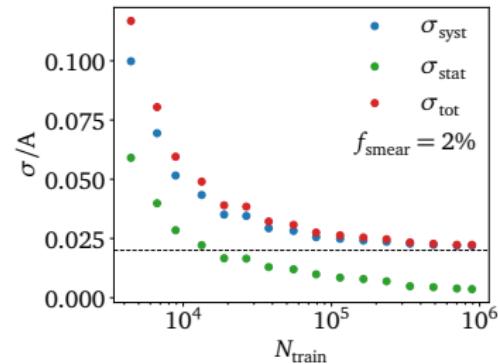
→ Another alternative for learned uncertainty



Amplitudes with calibrated uncertainties (A2b)

Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: artificial noise
- statistics plateau

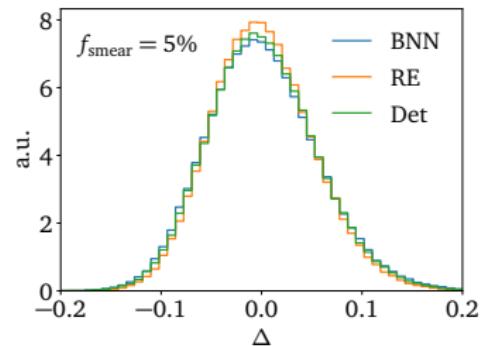
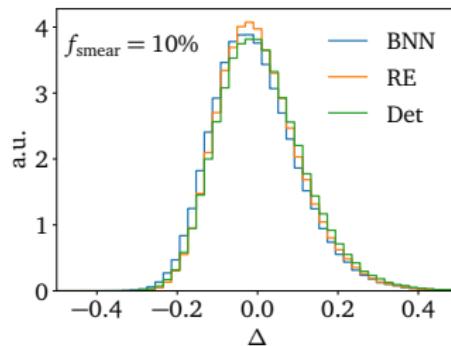


Amplitudes with calibrated uncertainties (A2b)

Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: artificial noise
- statistics plateau
- accuracy over phase space

$$\Delta(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{A_{\text{true}}(x)}$$



Amplitudes with calibrated uncertainties (A2b)

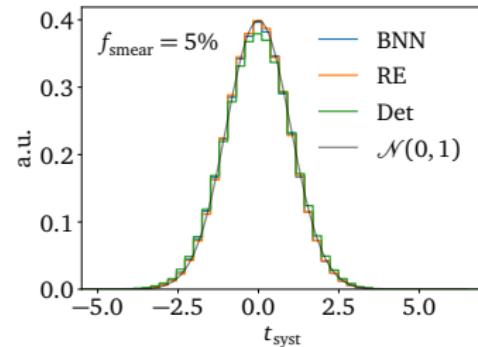
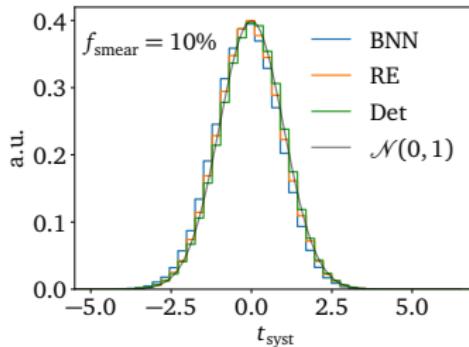
Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: **artificial noise**
- statistics plateau
- accuracy over phase space

$$\Delta(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{A_{\text{true}}(x)}$$

- pull over phase space

$$t_{\text{syst}}(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{\sigma_{\text{syst}}(x)}$$



Amplitudes with calibrated uncertainties (A2b)

Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: artificial noise
- statistics plateau
- accuracy over phase space

$$\Delta(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{A_{\text{true}}(x)}$$

- pull over phase space

$$t_{\text{syst}}(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{\sigma_{\text{syst}}(x)}$$

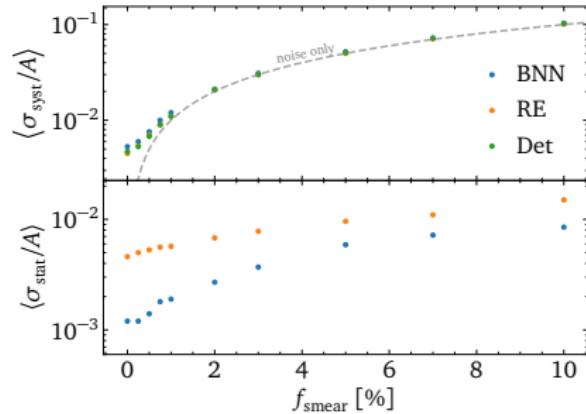
Towards zero noise

- scaling

$$\sigma_{\text{syst}}^2 - \sigma_{\text{syst},0}^2 \approx \sigma_{\text{train}}^2$$

- plateau $\langle \sigma_{\text{syst}} / A \rangle \sim 0.4\%$

→ Limiting factor??



Amplitudes with calibrated uncertainties (A2b)

Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: **artificial noise**
- statistics plateau
- accuracy over phase space

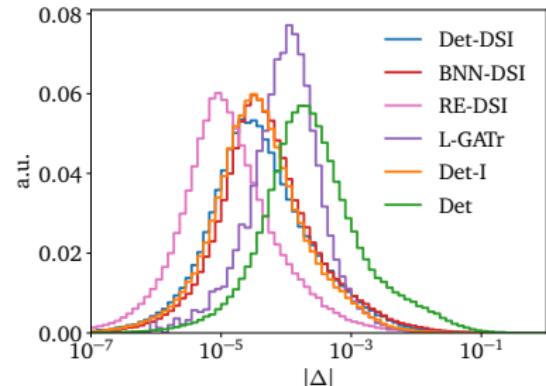
$$\Delta(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{A_{\text{true}}(x)}$$

- pull over phase space

$$t_{\text{syst}}(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{\sigma_{\text{syst}}(x)}$$

Data representation [Breso, Heinrich, Magerya, Olsson]

- amplitude from invariants
- learn Minkowski metric
- Deep-sets-invariant network



Amplitudes with calibrated uncertainties (A2b)

Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: **artificial noise**
- statistics plateau
- accuracy over phase space

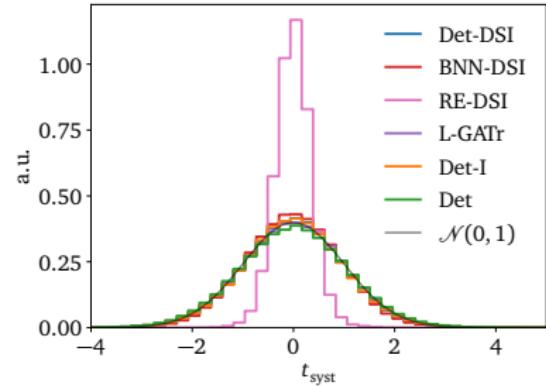
$$\Delta(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{A_{\text{true}}(x)}$$

- pull over phase space

$$t_{\text{syst}}(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{\sigma_{\text{syst}}(x)}$$

Data representation [Bresc, Heinrich, Magerya, Olsson]

- amplitude from invariants
 - learn Minkowski metric
 - Deep-sets-invariant network
- **Calibrated systematics**



Amplitudes with calibrated uncertainties (A2b)

Loop amplitude $gg \rightarrow \gamma\gamma g(g)$ [Bahl, Elmer, Favaro, Haußmann, TP, Winterhalder]

- systematics: **artificial noise**
- statistics plateau
- accuracy over phase space

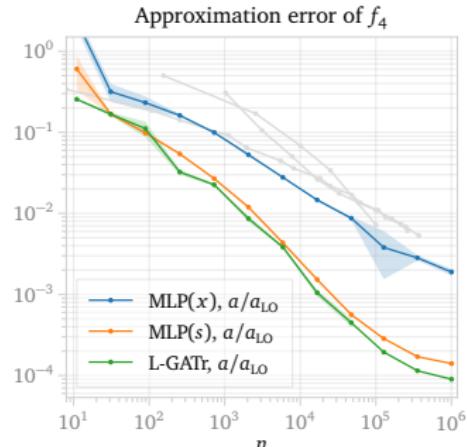
$$\Delta(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{A_{\text{true}}(x)}$$

- pull over phase space

$$t_{\text{syst}}(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{\sigma_{\text{syst}}(x)}$$

Data representation [Breso, Heinrich, Magerya, Olsson]

- amplitude from invariants
 - learn Minkowski metric
 - Deep-sets-invariant network
- **On to real loop integrals**



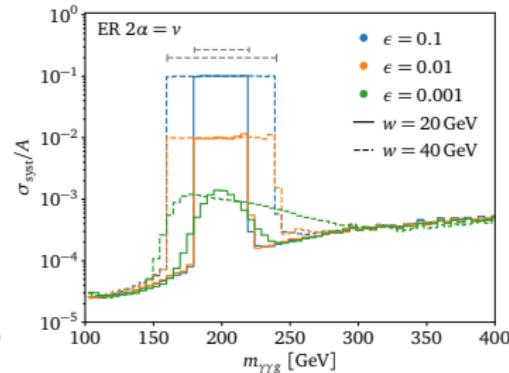
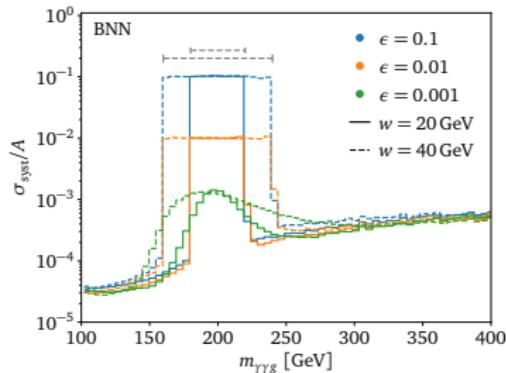
Phase space gaps

Noisy amplitude gap [Bahl, Elmer, TP, Winterhalder]

- step-function relative noise in $m_{\gamma\gamma g}$ gap

$$A_{\text{train}}(x) = \begin{cases} \mathcal{N}(A_{\text{true}}(x), \epsilon A_{\text{true}}(x)) & |m_{\gamma\gamma g}(x) - m_{\text{thresh}}| < w \\ A_{\text{true}} & |m_{\gamma\gamma g}(x) - m_{\text{thresh}}| \geq w \end{cases}$$

- compare Bayesian NN, ensembles, evidential regression
- Only little noise corrected, any noise learned



GANplification

Amplification from generative networks [Bahl, Diefenbacher, Elmer, TP, Spinner]

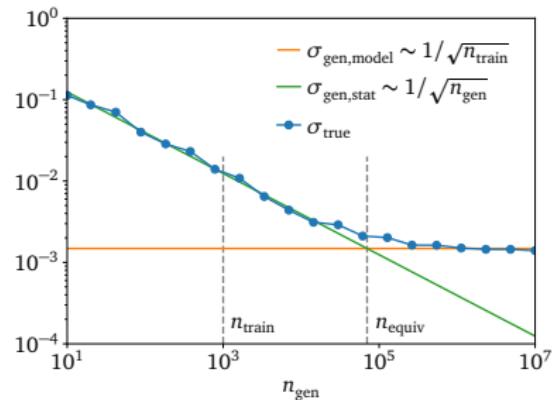
- phase space densities

$$p_{\text{gen}}(x) \approx p_{\text{train}}(x) \sim p_{\text{true}}(x) \quad \Rightarrow \quad p_{\text{gen}}(x) \stackrel{?}{\sim} p_{\text{true}}(x)$$

- 'How many events can I sample from a network trained on n events?'
fewer, unless trained perfectly, $p_{\text{gen}}(x) \neq p_{\text{train}}(x)$
more, generative network smoothen $p_{\text{train}}(x)$
- scaling of two uncertainties

$$\sigma_{\text{gen,stat}}^2 + \sigma_{\text{gen,model}}^2 = \begin{cases} \frac{a}{n_{\text{gen}}} & n_{\text{gen}} \ll n_{\text{train}} \\ \frac{b}{n_{\text{train}}} & n_{\text{gen}} \gg n_{\text{train}} \end{cases}$$

transition point $G \equiv \frac{n_{\text{gen}}}{n_{\text{train}}} = \frac{a}{b}$



GANplification

Amplification from generative networks [Bahl, Diefenbacher, Elmer, TP, Spinner]

- phase space densities

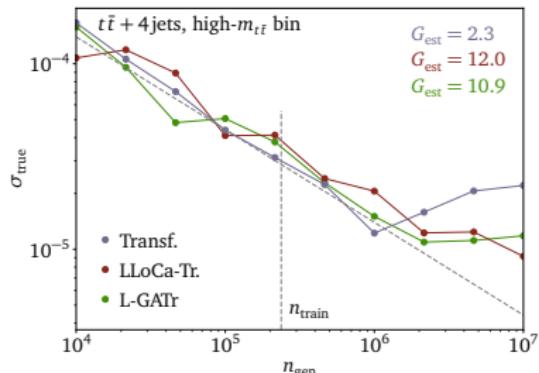
$$p_{\text{gen}}(x) \approx p_{\text{train}}(x) \sim p_{\text{true}}(x) \quad \Rightarrow \quad p_{\text{gen}}(x) \stackrel{?}{\sim} p_{\text{true}}(x)$$

- 'How many events can I sample from a network trained on n events?'
fewer, unless trained perfectly, $p_{\text{gen}}(x) \neq p_{\text{train}}(x)$
more, generative network smoothen $p_{\text{train}}(x)$
- scaling of two uncertainties

$$\sigma_{\text{gen,stat}}^2 + \sigma_{\text{gen,model}}^2 = \begin{cases} \frac{a}{n_{\text{gen}}} & n_{\text{gen}} \ll n_{\text{train}} \\ \frac{b}{n_{\text{train}}} & n_{\text{gen}} \gg n_{\text{train}} \end{cases}$$

transition point $G \equiv \frac{n_{\text{gen}}}{n_{\text{train}}} = \frac{a}{b}$

- averaging integral test
local KS-test
 - gen-equivalent training size
- L-GATr does amplify...



Physics from latent representation (B3b)

Quarks vs gluons from trained ParticleNet [Vent, Winterhalder, TP]

- sensitive substructure variables

$$n_{\text{pf}} = \sum_i 1 \quad w_{\text{pf}} = \frac{\sum_i p_{T,i} \Delta R_{i,\text{jet}}}{p_{T,\text{jet}}} \quad p_T D = \frac{\sqrt{\sum_i p_{T,i}^2}}{\sum_i p_{T,i}} \quad C_\beta = \frac{\sum_{i < j} p_{T,i} p_{T,j} (\Delta R_{ij})^\beta}{(\sum_i p_{T,i})^2}$$

- PC₁: constituent number and diversity

$$n_{\text{pf}} + \alpha \cdot S_{\text{PID}} \quad \text{with} \quad S_{\text{PID}} = - \sum_{\text{type } j} f_j \log f_j$$



Physics from latent representation (B3b)

Quarks vs gluons from trained ParticleNet [Vent, Winterhalder, TP]

- sensitive substructure variables

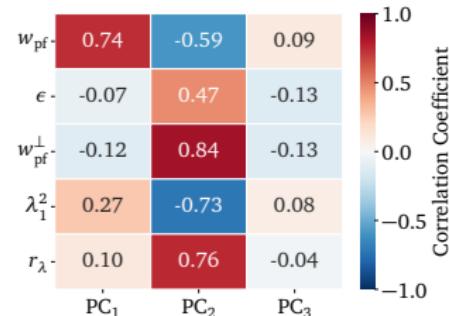
$$n_{\text{pf}} = \sum_i 1 \quad w_{\text{pf}} = \frac{\sum_i p_{T,i} \Delta R_{i,\text{jet}}}{p_{T,\text{jet}}} \quad p_T D = \frac{\sqrt{\sum_i p_{T,i}^2}}{\sum_i p_{T,i}} \quad C_\beta = \frac{\sum_{i < j} p_{T,i} p_{T,j} (\Delta R_{ij})^\beta}{(\sum_i p_{T,i})^2}$$

- PC₁: constituent number and diversity

$$n_{\text{pf}} + \alpha \cdot S_{\text{PID}} \quad \text{with} \quad S_{\text{PID}} = - \sum_{\text{type } j} f_j \log f_j$$

- PC₂: radial energy profile

$$w_{\text{pf}}^\perp = \alpha \cdot n_{\text{pf}} - w_{\text{pf}} \quad \text{and} \quad r_\lambda = \frac{\lambda_0^1}{\lambda_1^2} \quad \lambda_k^\beta = \sum_i z_i^\beta \Delta R^k$$



Physics from latent representation (B3b)

Quarks vs gluons from trained ParticleNet [Vent, Winterhalder, TP]

- sensitive substructure variables

$$n_{\text{pf}} = \sum_i 1 \quad w_{\text{pf}} = \frac{\sum_i p_{T,i} \Delta R_{i,\text{jet}}}{p_{T,\text{jet}}} \quad p_T D = \frac{\sqrt{\sum_i p_{T,i}^2}}{\sum_i p_{T,i}} \quad C_\beta = \frac{\sum_{i < j} p_{T,i} p_{T,j} (\Delta R_{ij})^\beta}{(\sum_i p_{T,i})^2}$$

- PC₁: constituent number and diversity

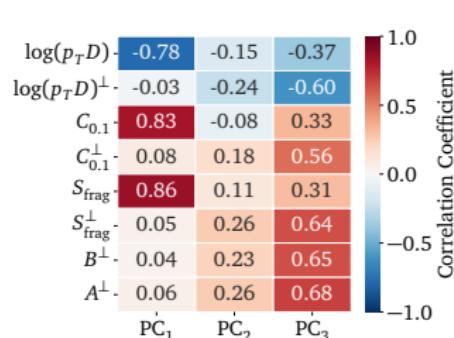
$$n_{\text{pf}} + \alpha \cdot S_{\text{PID}} \quad \text{with} \quad S_{\text{PID}} = - \sum_{\text{type } j} f_j \log f_j$$

- PC₂: radial energy profile

$$w_{\text{pf}}^\perp = \alpha \cdot n_{\text{pf}} - w_{\text{pf}} \quad \text{and} \quad r_\lambda = \frac{\lambda_0^1}{\lambda_1^2} \quad \lambda_k^\beta = \sum_i z_i^\beta \Delta R^k$$

- PC₃: fragmentation and energy dispersion

$$S_{\text{frag}} = - \sum_i z_i \log z_i$$



Physics from latent representation (B3b)

Symmetries

Uncertainties

Explainability

Anomalies

Quarks vs gluons from trained ParticleNet [Vent, Winterhalder, TP]

- sensitive substructure variables

$$n_{\text{pf}} = \sum_i 1 \quad w_{\text{pf}} = \frac{\sum_i p_{T,i} \Delta R_{i,\text{jet}}}{p_{T,\text{jet}}} \quad p_T D = \frac{\sqrt{\sum_i p_{T,i}^2}}{\sum_i p_{T,i}} \quad C_\beta = \frac{\sum_{i < j} p_{T,i} p_{T,j} (\Delta R_{ij})^\beta}{(\sum_i p_{T,i})^2}$$

- PC₁: constituent number and diversity

$$n_{\text{pf}} + \alpha \cdot S_{\text{PID}} \quad \text{with} \quad S_{\text{PID}} = - \sum_{\text{type } j} f_j \log f_j$$

- PC₂: radial energy profile

$$w_{\text{pf}}^\perp = \alpha \cdot n_{\text{pf}} - w_{\text{pf}} \quad \text{and} \quad r_\lambda = \frac{\lambda_{0.5}^1}{\lambda_1^2} \quad \lambda_k^\beta = \sum_i z_i^\beta \Delta R^k$$

- PC₃: fragmentation and energy dispersion

$$S_{\text{frag}} = - \sum_i z_i \log z_i$$

- PC_{4,5}: charge information etc

$$E_Q = \frac{E_{\text{charged}}}{E_{\text{jet}}} \quad \text{and} \quad A^\perp = S_{\text{frag}} \frac{C_{0.1}}{C_{0.05}} - 0.03 \cdot n_{\text{pf}} + 1.95 w_{\text{pf}}^\perp$$

→ Latent distributions learn physics



ParticleNet beyond PCA

Disentangled latent classifier

- learning compressed, decorrelated representation

$$\mathcal{L} = \underbrace{\sum_{i=1}^N |x_i - \hat{x}_i|^2}_{\mathcal{L}_{\text{reco}}} + \underbrace{\sum_{i=1}^N [y_i \log \sigma(z_i) + (1 - y_i) \log(1 - \sigma(z_i))]}_{\mathcal{L}_{\text{class}}} + \underbrace{\sum_{j \neq k} [\text{Cov}(z_j, z_k)]^2}_{\mathcal{L}_{\text{disentangle}}}$$

→ 5 latent dimensions plenty

Latent Dim	1	2	3	4
AUC	0.893(2)	0.9001(4)	0.9024(4)	0.9034(2)
rej _{30%}	72(3)	77(3)	95(5)	95(3)
ΔC	1.8(3)	0.93(5)	1.0(16)	0.9(15)

Symbolic Regression

- learn formula with given complexity
- classifier output not power series

→ Formulas as physics regularizers? [Bahl, Fuchs, Menem, TP]

$$p_{\text{quark}} = \tanh^3 \left[0.55 \cdot C_{0.2} + 2 \left(-0.02 \cdot r_\lambda \cdot (C_{0.2} \cdot p_T D \cdot S_{\text{PID}} \cdot S_{\text{frag}} - 0.25) + 1 \right)^3 \right]$$

observables	model	AUC	Rej _{30%}
(n_{pf} , $p_T D$, $C_{0.2}$, r_λ , S_{PID} , S_{frag} , E_Q)	MLP PySR	0.872 0.871	66.87 66.58



Extrapolating ISR

Universal QCD jet radiation [Z + 1...8 jets]

- from n to $n+1$ jets

$$R_{(n+1)/n} = \frac{\sigma_{n+1}}{\sigma_n} \quad \text{and} \quad P(n) = \frac{\sigma_n}{\sigma_{\text{tot}}} \quad \text{with} \quad \sigma_{\text{tot}} = \sum_{n=0}^{\infty} \sigma_n .$$

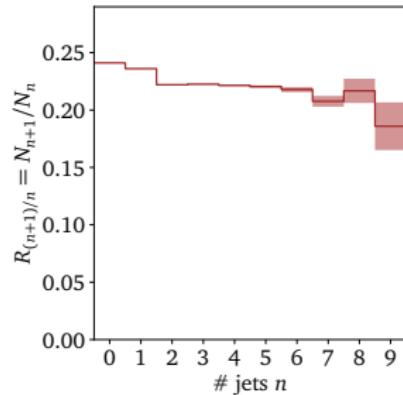
- large scale drop: Poisson scaling

$$R_{(n+1)/n} = \frac{\bar{n}}{n+1} \quad \Leftrightarrow \quad P(n) = \frac{\bar{n}^n e^{-\bar{n}}}{n!} .$$

- democratic scales: staircase scaling

$$R_{(n+1)/n} = e^{-b} = 1 - \tilde{\Delta}_g(Q^2) \equiv P(n+1|n)$$

→ Universal pattern learnable?



Extrapolating ISR

Universal QCD jet radiation $[Z + 1 \dots 8 \text{ jets}]$

- democratic scales: staircase scaling

$$R_{(n+1)/n} = e^{-b} = 1 - \tilde{\Delta}_g(Q^2) \equiv P(n+1|n)$$

→ Universal pattern learnable?

Autoregressive transformer $[$ Butter, Charton, Villadamigo, Ore, TP, Spinner $]$

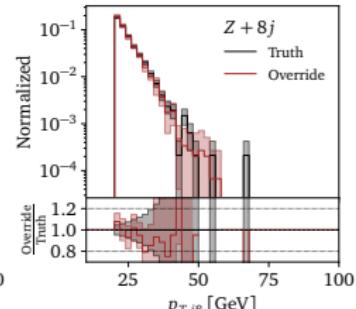
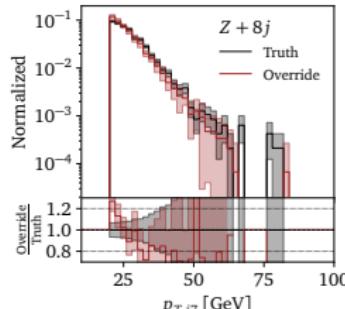
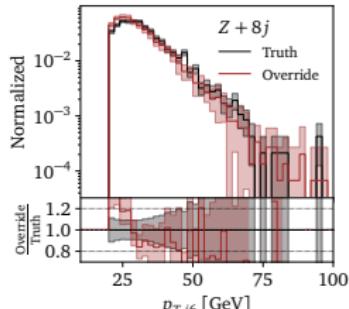
- factorized probability and loss function

$$p(x_i|x_{1:i-1}) = p_{\text{kin}}(x_i|x_{1:i-1}) p_{\text{split}}(x_{1:i-1})$$

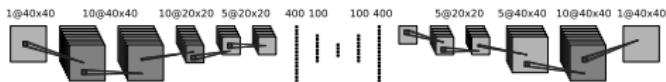
$$p(x_{1:n}) = \left[\prod_{i=1}^n p_{\text{kin}}(x_i|x_{1:i-1}) \right] \left[\prod_{i=1}^n p_{\text{split}}(x_{1:i-1}) \right] [1 - p_{\text{split}}(x_{1:n})] ,$$

- train up to 6 jets, generate 7 and 8 jets

→ full extrapolation with right latent representation



Anomaly searches (B3b)



Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets → top jets hard to describe
 - reconstruct top jets → QCD jets just simple top-like jet
- Symmetric performance $S \leftrightarrow B?$



Anomaly searches (B3b)

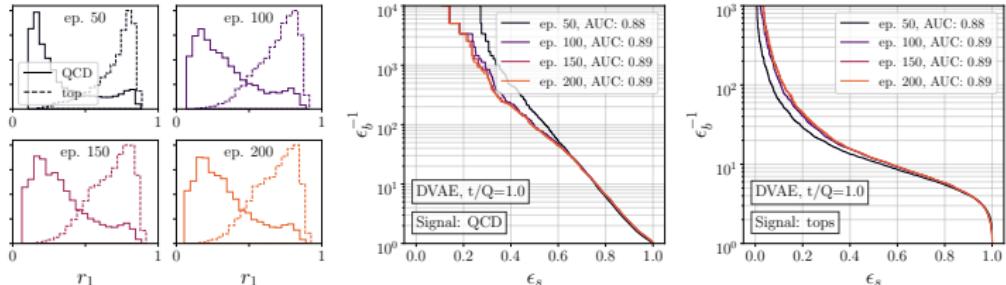


Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets → top jets hard to describe
 - reconstruct top jets → QCD jets just simple top-like jet
- Symmetric performance $S \leftrightarrow B$?

Moving to latent space [Dillon, Favaro, TP, Sorrensen, Krämer]

- anomaly score from latent space?
- VAE → does not work
Gaussian mixture VAE → does not work
Dirichlet VAE → works okay
density estimation → does not work



Anomaly searches (B3b)

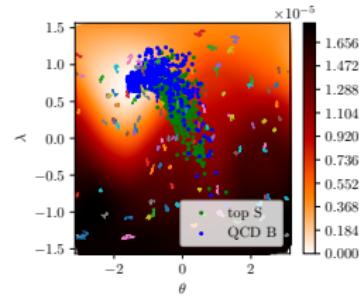
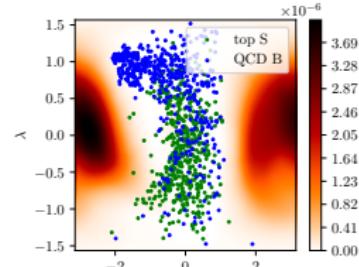


Unsupervised classification

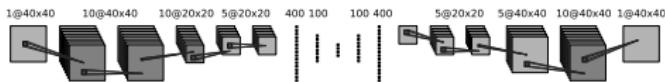
- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets → top jets hard to describe
 - reconstruct top jets → QCD jets just simple top-like jet
- Symmetric performance $S \leftrightarrow B$?

Normalized autoencoder [Sangwoong Yoon, Noh, Park]

- normalized probability loss
 - Boltzmann mapping [$E_\theta = \text{MSE}$]
- $$p_\theta(x) = \frac{e^{-E_\theta(x)}}{Z_\theta}$$
- $$L = -\langle \log p_\theta(x) \rangle = \langle E_\theta(x) + \log Z_\theta \rangle$$
- inducing background metric
 - large MSE for too much and missing structure
- Symmetric autoencoder, deployed by CMS



Anomaly searches (B3b)



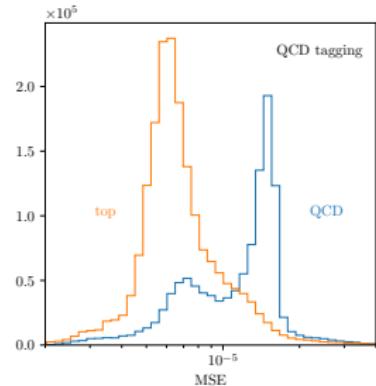
Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets → top jets hard to describe
 - reconstruct top jets → QCD jets just simple top-like jet
- Symmetric performance $S \leftrightarrow B?$

Normalized autoencoder [Sangwoong Yoon, Noh, Park]

- normalized probability loss
 - Boltzmann mapping [$E_\theta = \text{MSE}$]
$$p_\theta(x) = \frac{e^{-E_\theta(x)}}{Z_\theta}$$

$$L = -\langle \log p_\theta(x) \rangle = \langle E_\theta(x) + \log Z_\theta \rangle$$
 - inducing background metric
 - large MSE for too much and missing structure
- Symmetric autoencoder, deployed by CMS



AI for fundamental physics

Develop AI for the best science

1 just another tool for a numerical field

2 transformative new language

- many applications in LHC theory

MadGraph7

MLhad

Higher orders (A2b)

Simulation-based inference (A3a)

SFitter global analyses (A2a)

Unfolding

...

→ Make complexity our friend

Modern Machine Learning for LHC Physicists

Tilman Plehn^{a*}, Anja Butter^{a,b}, Barry Dillon^a,
Theo Heimel^a, Claudius Krause^c, and Ramon Winterhalder^d

^a Institut für Theoretische Physik, Universität Heidelberg, Germany

^b LPNHE, Sorbonne Université, Université Paris Cité, CNRS/IN2P3, Paris, France

^c HEPHY, Austrian Academy of Sciences, Vienna, Austria

^d CP3, Université catholique de Louvain, Louvain-la-Neuve, Belgium

March 19, 2024

Abstract

Modern machine learning is transforming particle physics fast, bullying its way into our numerical tool box. For young researchers it is crucial to stay on top of this development, which means applying cutting-edge methods and tools to the full range of LHC physics problems. These lecture notes lead students with basic knowledge of particle physics and significant enthusiasm for machine learning to relevant applications. They start with an LHC-specific motivation and a non-standard introduction to neural networks and then cover classification, unsupervised classification, generative networks, and inverse problems. Two themes defining much of the discussion are well-defined loss functions and uncertainty-aware networks. As part of the applications, the notes include some aspects of theoretical LHC physics. All examples are chosen from particle physics publications of the last few years.¹



AI for fundamental physics

Develop AI for the best science

- 1 just another tool for a numerical field
- 2 transformative new language

- many applications in LHC theory

MadGraph7

MLhad

Higher orders (A2b)

Simulation-based inference (A3a)

SFitter global analyses (A2a)

Unfolding

...

→ And find a new job...

Agents of Discovery

Sascha Diefenbacher¹, Anna Hallin²,
Gregor Kasieczka², Michael Krämer³, Anne Lauscher⁴, Tim Lukas²,

¹ Physics Division, Lawrence Berkeley National Laboratory, Berkeley, USA
² Institut für Experimentalphysik, Universität Hamburg, Germany

³ Institute for Theoretical Particle Physics and Cosmology, RWTH Aachen University,
Germany

⁴ Data Science Group, Universität Hamburg, Germany

September 11, 2025

Abstract

The substantial data volumes encountered in modern particle physics and other domains of fundamental physics research allow (and require) the use of increasingly complex data analysis tools and workflows. While the use of machine learning (ML) tools for data analysis has recently proliferated, these tools are typically special-purpose algorithms that rely, for example, on encoded physics knowledge to reach optimal performance. In this work, we investigate a new and orthogonal direction: Using recent progress in large language models (LLMs) to create a team of *agents* — instances of LLMs with specific subtasks — that jointly solve data analysis-based research problems in a way similar to how a human researcher might: by creating code to operate standard tools and libraries (including ML systems) and by building on results of previous iterations. If successful, such agent-based systems could be deployed to automate routine analysis components to counteract the increasing complexity of modern tool chains. To investigate the capabilities of current-generation commercial LLMs, we consider the task of anomaly detection via the publicly available and highly-studied LHC Olympics dataset. Several current models by OpenAI (GPT-4o, o4-mini, GPT-4.1, and GPT-5) are investigated and their stability tested. Overall, we observe the capacity of the agent-based system to solve this data analysis problem. The best agent-created solutions mirror the performance of human state-of-the-art results.

