



Personel

- In Y2: Sebastian (WBS1.4.5) was $\frac{1}{3}$ FTE on IRIS-HEP IA (other $\frac{2}{3}$ on DIANA-HEP)
- Heiko & Irina are funded via SCALFIN, some projects interact with IRIS-HEP (eg. Reproducible Open Benchmarks (ROB), REANA workflows for MadMiner)
- Sinclert (WBS1.5.5) is $\sim\frac{1}{3}$ FTE on IRIS-HEP, currently contributing to AS/SSC, but possible contributions to ROB or something similar.



Kyle Cranmer
New York University



**Sebastian
Macaluso**
New York University



Heiko Mueller
New York University

*Research Software
Engineer (SCALFIN)*



Sinclert Pérez
New York University

*Research Software
Engineer*



Irina Espejo
New York University



Achievements



Machine Learning for Jet Physics Workshop

Sebastian Macaluso and I developed:

- **Ginkgo**, a simplified parton shower developed in a probabilistic programming language to streamline research with CS / stats / ML community
- **Hierarchical Trellis** algorithm for efficiently searching and summing over binary trees
- Contributions to Reproducible Open Benchmarks (ROB)

Ginkgo: Toy Generative Model for Jets



Generative model to aid in ML research for jet physics.

NLP analogy: ground-truth parse trees with a known language model

K. Cranmer, S. Macaluso & D. Pappadopulo

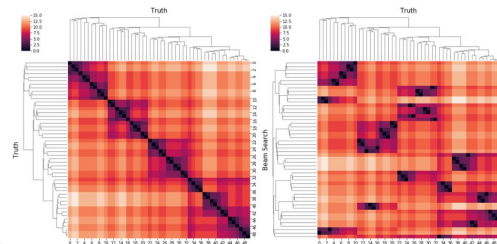
github.com/SebastianMacaluso/ToyJetsShower

Greedy Algorithm

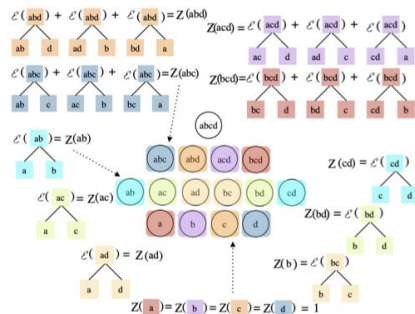
Locally maximizing the likelihood at each step.

Beam Search Algorithm

Maximize the likelihood of multiple steps before choosing the latent path.



Use-inspired Research: Hierarchical Cluster Trellis for Exact Inference



Standard clustering algorithm in HEP (anti-kT) is Greedy

$$P_{\text{Greedy}} < P_{\text{Beam Search}} < P_{\text{Trellis}}$$

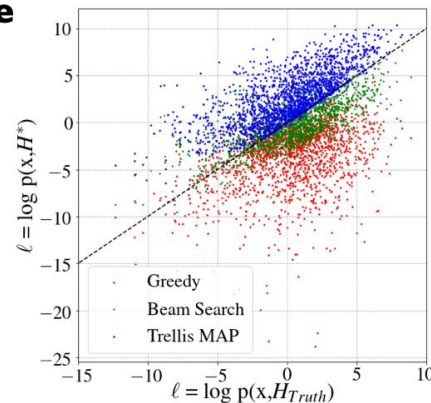
Use inspired-research

Seminar at UMass Amherst, Center for Data Science, College of Information & Computer Sciences led to collaboration for hierarchical clustering algorithms.

S. Macaluso, C. Greenberg, N. Monath, J. Lee, P. Flaherty, K. Cranmer, A. McGregor, A. McCallum

Applications in other domains, e.g. cancer genomics.

<https://arxiv.org/abs/2002.11661>





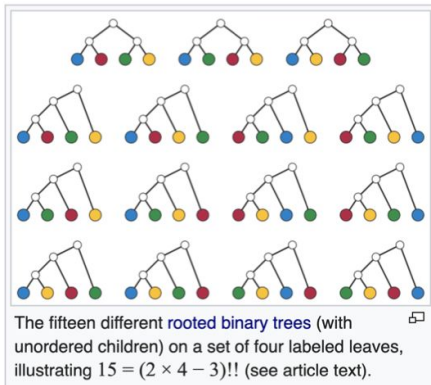
A word on hierarchical trellis

- Many tasks (eg. jet clustering) include trying to find the best hierarchical clustering over N objects, or summing over all of them weighted in some way (eg. CKKW-L matrix element / parton shower matching, and matrix element method)
- Number of trees grows like a double-factorial !!
- Our trellis algorithm allows us to do max/sum efficiently and sparse trellis can scale.

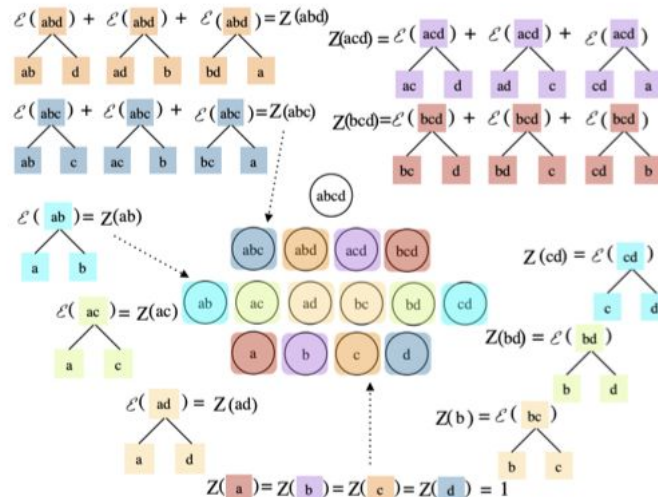
Challenge

Number of clustering histories
for N leaves grows as
 $a(N) = (2N - 3)!!$

# of leaves	Approx. # of trees
4	15
5	100
7	10 k
9	2 M
11	600 M



https://en.wikipedia.org/wiki/Double_factorial





Milestones & Metrics

Status of Year 1 and 2 milestones:

- G4.12: Evaluation of NN exchange formats (eg. ONNX, torchscript)
 - Both ATLAS and CMS have advanced on this
 - Was discussed in FastML blueprint.
 - Need to pound out a short document to summarize

Potential Y3 Milestones (remember small effort):

- Performance of sparse hierarchical trellis algorithm (compared to greedy, beam search, exact trellis in terms of finding maximum likelihood clustering and marginal likelihood)
- Snowmass white paper contributions about potential to unify Monte Carlo generation (eg. parton shower) and inference (eg. jet tagging) with probabilistic programming
- Exploratory:
 - Exploration of hierarchical trellis to speed up the bottleneck in the CKKW-L matrix element / parton shower model (investigation of generalization to 2->3 splittings)
 - Exploration of probabilistic programming for MC generation in tails of jet phase space

Input on Metrics: would be good to have metrics that pair well with more exploratory work & efforts to show impact beyond HEP [papers, external collaborations, cross-over].





Challenges, Outcomes, Feedback

Ideas for Grand Challenges:

- Could potentially include hierarchical trellis jet reconstruction algorithms or some ML tagging algs in a reconstruction-focused grand challenge

Proposed project outcomes:

- Roadmap / white paper outlining what new capabilities would be enabled by enabling Monte Carlo generators (focusing on Jets) with probabilistic programming functionality and what interventions in code are needed to get there:
 - Improved MC tuning, remove CKKW-L bottle neck, more efficient MC generation in tails of jet substructure, better jet clustering and tagging, etc.

Feedback:

- AS has biweekly meetings, which is nice for coherence, and we have a postdoc chair the meeting, which is nice for professional development.
 - Not clear that carries over well to IA, but food for thought.

