

Identification of Jets Containing b -Hadrons with Recurrent Neural Networks at the ATLAS Experiment

ATL-PHYS-PUB-2017-003

Dan Guest
ATLAS Collaboration

UC Irvine

May 9, 2017

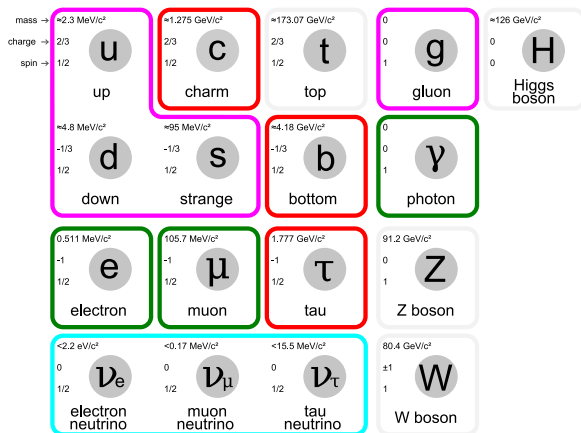
Why b -tag? (an oversimplification)

- ▶ Want Higgs?
- ▶ **Higgs** mostly decays to **b -quarks**
- ▶ b -quarks make jets
- ▶ In LHC, *everything* makes jets
- ▶ Not everything makes b -jets

mass →	$\approx 2.3 \text{ MeV}/c^2$	$\approx 1.275 \text{ GeV}/c^2$	$\approx 173.07 \text{ GeV}/c^2$	0	$\approx 126 \text{ GeV}/c^2$
charge →	$2/3$	$2/3$	$2/3$	0	0
spin →	$1/2$	$1/2$	$1/2$	1	0
	u up	c charm	t top	g gluon	H Higgs boson
QUARKS					
	$\approx 4.8 \text{ MeV}/c^2$	$\approx 95 \text{ MeV}/c^2$	$\approx 4.18 \text{ GeV}/c^2$	0	
	$-1/3$	$-1/3$	$-1/3$	0	
	$1/2$	$1/2$	$1/2$	1	
	d down	s strange	b bottom	γ photon	
	$0.511 \text{ MeV}/c^2$	$105.7 \text{ MeV}/c^2$	$1.777 \text{ GeV}/c^2$	$91.2 \text{ GeV}/c^2$	
	-1	-1	-1	0	
	$1/2$	$1/2$	$1/2$	1	
	e electron	μ muon	τ tau	Z Z boson	
	$< 2.2 \text{ eV}/c^2$	$< 0.17 \text{ MeV}/c^2$	$< 15.5 \text{ MeV}/c^2$	$80.4 \text{ GeV}/c^2$	
	0	0	0	± 1	
	$1/2$	$1/2$	$1/2$	1	
	ν_e electron neutrino	ν_μ muon neutrino	ν_τ tau neutrino	W W boson	
LEPTONS				GAUGE BOSONS	

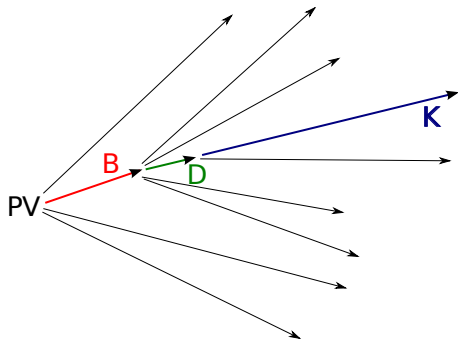
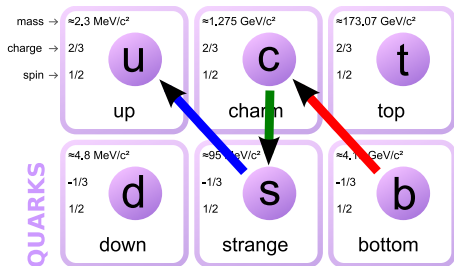
The Standard Model (as Seen by Collider Physics)

- ▶ **Some are stable**
- ▶ Many unstable
- ▶ **Some form jets**
- ▶ **Some metastable**
- ▶ Neutrinos $\rightarrow E_T^{\text{miss}}$



- ▶ Short-lived particles are a big part of what we measure!

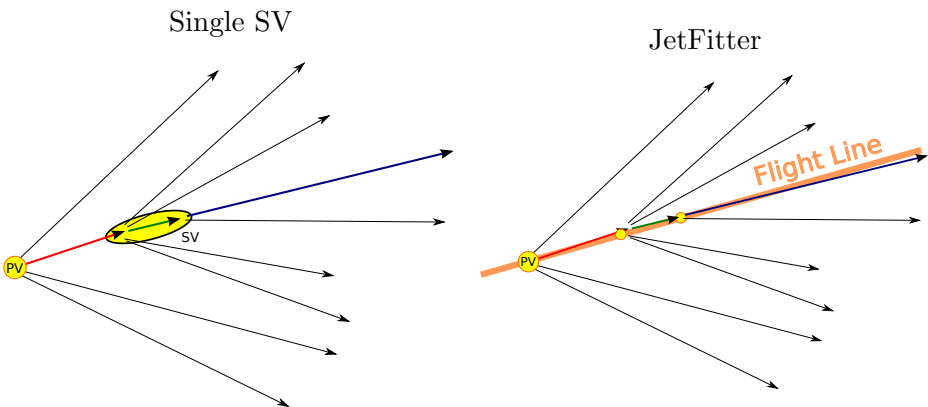
The b -hadron Decay Chain



- ▶ b -hadrons decay through cascade
- ▶ $\beta\gamma c\tau \approx 6.4 \text{ mm}$ for B with $p_T = 70 \text{ GeV}$
- ▶ But many decay distances are $O(\text{detector resolution})$

Reconstructing Secondary Vertices

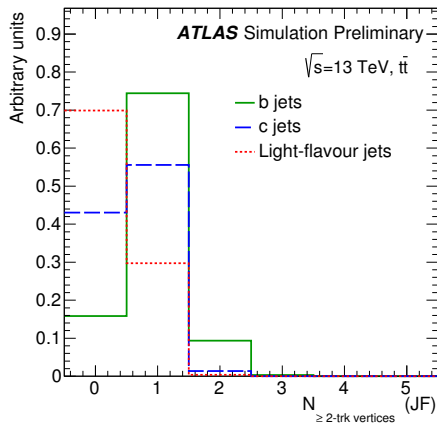
The ATLAS approaches



- ▶ Many discriminants come from vertices, combine them with ML

The problem with SV tagging

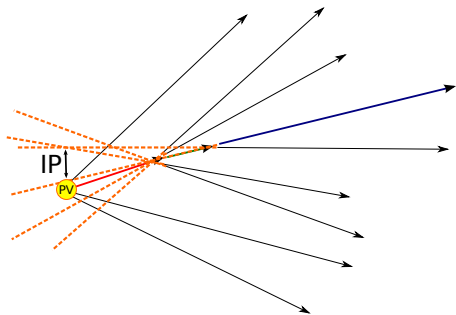
- ▶ Sometimes we don't find a vertex
- ▶ Requires cutting on track-vertex compatibility
 - ▶ Cuts always lose information
- ▶ Tuned “by hand”
- ▶ Experiment-specific



- ▶ There is no FASTJET for vertex reconstruction

Impact parameter (IP) tagging

- ▶ Take all tracks in a jet
- ▶ Apply some selection
- ▶ Extrapolate to perigee
- ▶ Per-track discriminants:
 - ▶ $S_{d_0} \equiv d_0/\sigma_{d_0}$
 - ▶ $S_{z_0} \equiv z_0/\sigma_{z_0}$
 - ▶ track “quality”
- ▶ Compute per-track likelihood $\mathcal{L}_f(\text{track})$ with $f \in \{b, c, \text{light}\}$
- ▶ Per-jet likelihood $p_f = \prod_{\text{trk}} \mathcal{L}_f(\text{track variables})$
- ▶ **IP based tagging is the problem we solve with RNNs**
 - ▶ More on this later



Putting it all together

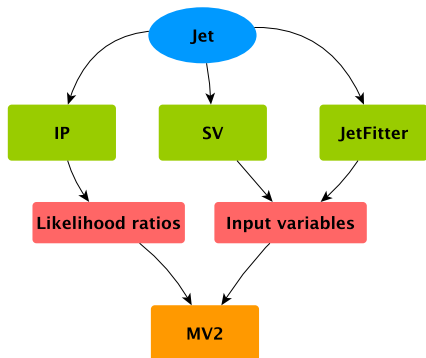
Low-Level

- ▶ **IP**: track-based variables
- ▶ **Likelihood**: gives $p_b, p_c, p_{\text{light}}$
- ▶ **SV**: gives vertex variables
- ▶ **JetFitter**: similar to SVx

High-level

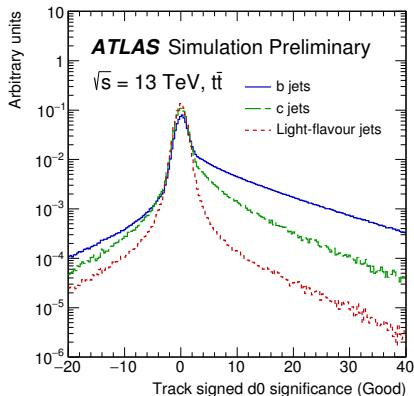
- ▶ **MV2**: combine with BDT

- ▶ It's easy to focus on the high-level tagger (MV2), but upstream is important too



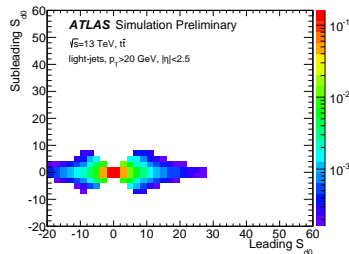
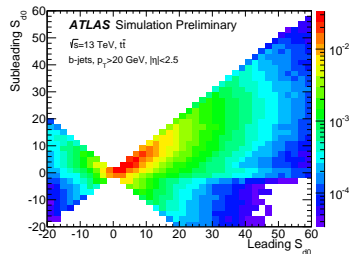
IP3D: ATLAS's IP Tagger

- ▶ Need to define $\mathcal{L}_f(\text{track})$
 - ▶ $\mathcal{L}_f(S_{d_0}, S_{z_0}, \text{category})$
 - ▶ S_{d_0} shown right
- ▶ Use histograms from simulation
- ▶ 3D binning scheme:
 - ▶ 35 bins in S_{d_0}
 - ▶ 20 bins in S_{z_0}
 - ▶ 14 track categories
- ▶ track category represents quality of track



Improving Upstream Taggers: What IP3D misses

- ▶ Relations among tracks:
 - ▶ relation to neighbor bins
 - ▶ relation to neighbor tracks
- ▶ **These are important** (see right)
- ▶ New (SV inspired) track variables:
 - ▶ $p_T^{\text{frac}} \equiv p_T^{\text{track}} / p_T^{\text{jet}}$
 - ▶ $\Delta R(\text{track}, \text{jet})$

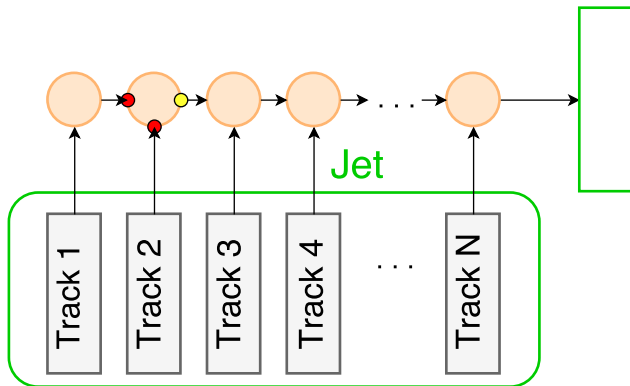


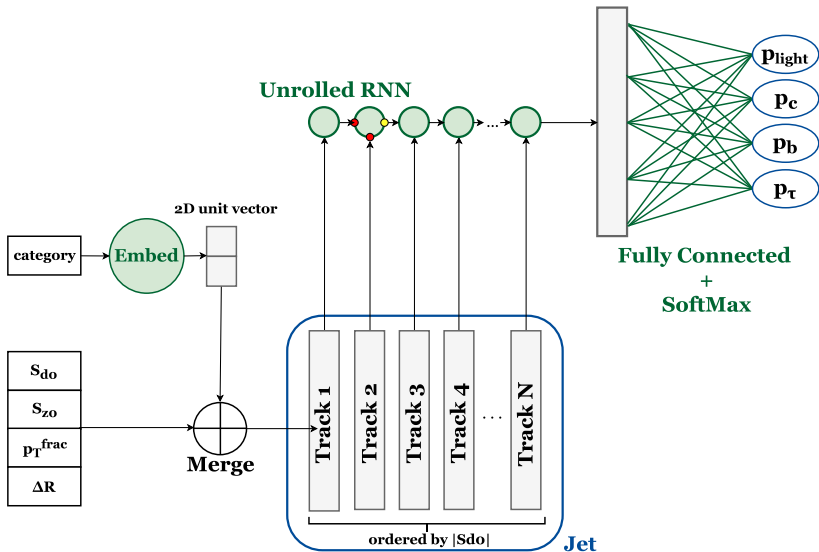
Curse of Dimensionality

- ▶ Already 29,400 bins
- ▶ New variable $\rightarrow \sim 10 \times$ bins (and events to “train”)

Recurrent Neural Networks (RNNs)

- ▶ RNNs can process an arbitrarily length sequence
- ▶ Output is a fixed dimensional vector for each jet





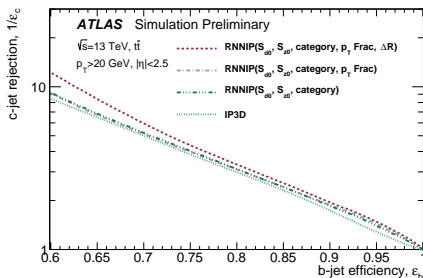
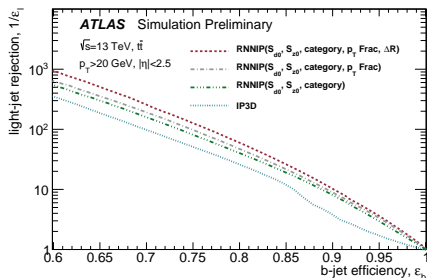
ROC Curves for a Multi-Background Discriminant

- ▶ Eventually we'll combine with vertex-based approaches
- ▶ Conventional HEP discriminants are binary
 - ▶ Train against a mix of backgrounds (i.e. MV2 is 7% c -jets)
- ▶ We use 4 outputs:
 - ▶ p_b : bottom jet
 - ▶ p_c : charm jet
 - ▶ p_{light} : “light” jet (u, d, s, g)
 - ▶ p_τ : τ jet
- ▶ Combine everything for the sake of plots

$$D_{\text{RNN}} = \ln \frac{p_b}{f_c p_c + f_\tau p_\tau + (1 - f_c - f_\tau) p_{\text{light}}} \quad (1)$$

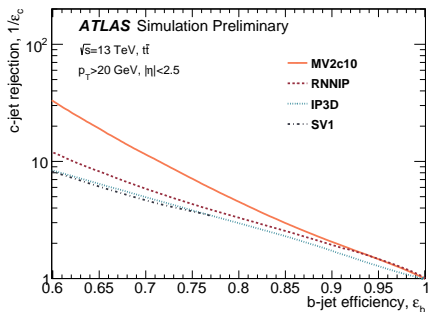
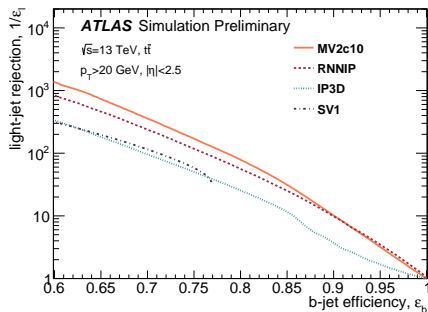
- ▶ The f weighting parameters can be adjusted post-training
- ▶ For this talk: $f_c = 0.07, f_\tau = 0$

RNN Performance (compared to IP3D)



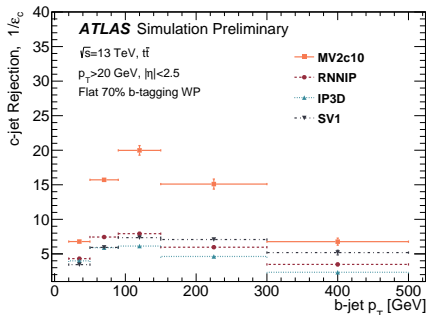
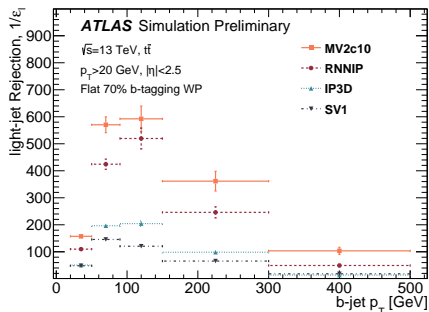
- ▶ Lowest line is IP3D
- ▶ Next up: RNN with IP3D inputs
- ▶ Each new variable adds discrimination
- ▶ At 70% working point:
 - ▶ RNN with IP3D inputs improves light rejection by 1.7
 - ▶ With $\Delta R(\text{track}, \text{jet})$ and p_T^{frac} , improves light rejection by 2.5

RNN Performance (compared to high-level tagger)



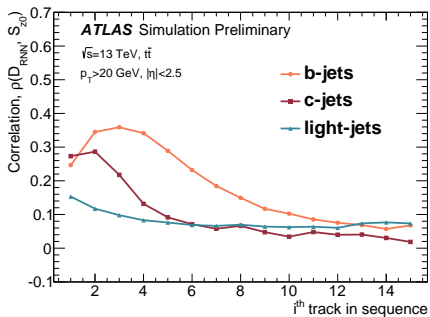
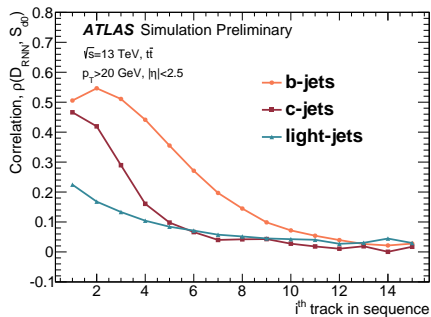
- ▶ MV2 using IP3D still rejects more background for $\epsilon_b < 0.9$
- ▶ But this uses JetFitter and SV \rightarrow much more information
- ▶ RNN as input for MV2 is outside the scope of this talk
 - ▶ But we can imagine replacing IP3D with the RNN

RNN Performance by p_T



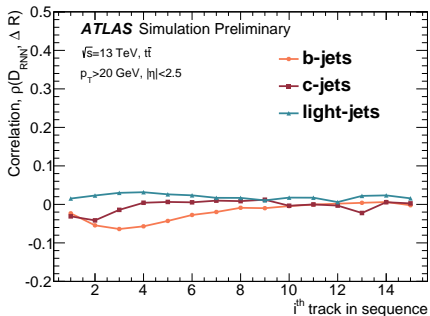
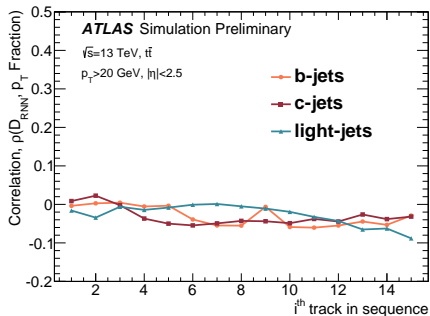
- ▶ Cut on the discriminant such that $\varepsilon_b = 0.7$ in each p_T bin
- ▶ Same trend as previous slide: rejection for IP3D < RNN < MV2
- ▶ RNN tagger is no more p_T dependent than other taggers

RNN output correlation with input: S_{d_0} and S_{z_0}



- ▶ D_{RNN} output is highly correlated with jet S_{d_0} for “early” tracks in $|S_{d_0}|$ ordering
 - ▶ Interesting, but maybe not surprising: b hadrons have ~ 5 tracks
- ▶ Effect is less pronounced for S_{z_0}

RNN output correlation with input: ΔR and p_T^{frac}



- ▶ Much less correlation between D_{RNN} and $\Delta R(\text{track}, \text{jet})$ or p_T^{frac}
- ▶ But these are useful discriminants nonetheless

Notes on Software

- ▶ We train with KERAS
 - ▶ Use THEANO backend
- ▶ **Our reconstruction framework doesn't support batched NumPy arrays**
- ▶ Within our reconstruction, we evaluate with LWTNN
 - ▶ Used in ATLAS for top and W tagging
 - ▶ Also used by CMS for DeepFlavour

Help and other ideas welcome

- ▶ LWTNN is written “as needed”
- ▶ Is there a more sustainable approach?

Conclusions

- ▶ RNNs are a promising tool for flavor tagging
 - ▶ Use relatively low-level variables
 - ▶ Can augment vertex-based approaches
- ▶ Many interesting questions:
 - ▶ What other low-level variables could we include?
 - ▶ How does this complement a high-level tagger (e.g. MV2, DeepFlavour)?
 - ▶ How does this compare to the CMS approach?
 - ▶ Can we “understand” (visualize) what we’ve learned?
- ▶ Thanks for listening, ideas are welcome!

BACKUP

Thanks

- ▶ Michela Paganini and Jonathan Shlomi for the graphics
- ▶ Zihao Jiang, Michael Kagan, Michela, and the rest of the RNN team for training lots of networks
- ▶ The ATLAS flavor tagging group for a good problem
- ▶ ATLAS for all the simulation

IP3D Categories

#	Category	Fractional contribution [%]		
		<i>b</i> -jets	<i>c</i> -jets	light-jets
0	No hits in first two layers; expected hit in IBL and b-layer	1.9	2.0	1.9
1	No hits in first two layers; expected hit in IBL and no expected hit in b-layer	0.1	0.1	0.1
2	No hits in first two layers; no expected hit in IBL and expected hit in b-layer	0.04	0.04	0.04
3	No hits in first two layers; no expected hit in IBL and b-layer	0.03	0.03	0.03
4	No hit in IBL; expected hit in IBL	2.4	2.3	2.1
5	No hit in IBL; no expected hit in IBL	1.0	1.0	0.9
6	No hit in b-layer; expected hit in b-layer	0.5	0.5	0.5
7	No hit in b-layer; no expected hit in b-layer	2.4	2.4	2.2
8	<i>Shared</i> hit in both IBL and b-layer	0.01	0.01	0.03
9	At least one <i>shared</i> pixel hits	2.0	1.7	1.5
10	Two or more <i>shared</i> SCT hits	3.2	3.0	2.7
11	<i>Split</i> hits in both IBL and b-layer	1.0	0.87	0.6
12	<i>Split</i> pixel hit	1.8	1.4	0.9
13	<i>Good</i>	83.6	84.8	86.4

- Fractions are based on simulated $t\bar{t}$

Training

- ▶ Use 3.2 million jets from simulated $t\bar{t}$
- ▶ Training time: with a CPU, a few days on a (busy) cluster
- ▶ We only train on first 15 tracks (0.5% of jets 15+ tracks)

Track Selection

- ▶ Jet Algorithm: Anti- k_t , $R = 0.4$
- ▶ Track $p_T > 1$ GeV
- ▶ $|d_0| < 1$ mm, $|z_0 \sin \theta| < 1.5$ mm
- ▶ $n_{\text{hits}}^{\text{Si}} \geq 7$, $n_{\text{holes}}^{\text{Si}} \leq 2$, $n_{\text{holes}}^{\text{pixel}} \leq 1$