

Lecture 4: A Case Study in Jet Substructure

The last decade has seen a transformation in how we think about jets and jet substructure. Many clever observables have been developed to improve experimental robustness/performance and enable theoretical control. There is a growing catalog of precise jet measurements/calculations, and with new machine learning techniques, no signs of a slow down.

I want to end these lectures with a classic but simple calculation that uses the ingredients you've learned. It is a jet discrimination task that has a history going back almost four decades.

Quark Jets vs. Gluon Jets.

At a cartoon level, this problem is simple.

$$q \begin{array}{c} \diagup \\ \diagdown \\ \diagup \\ \diagdown \end{array} \quad C_F = 4/3$$

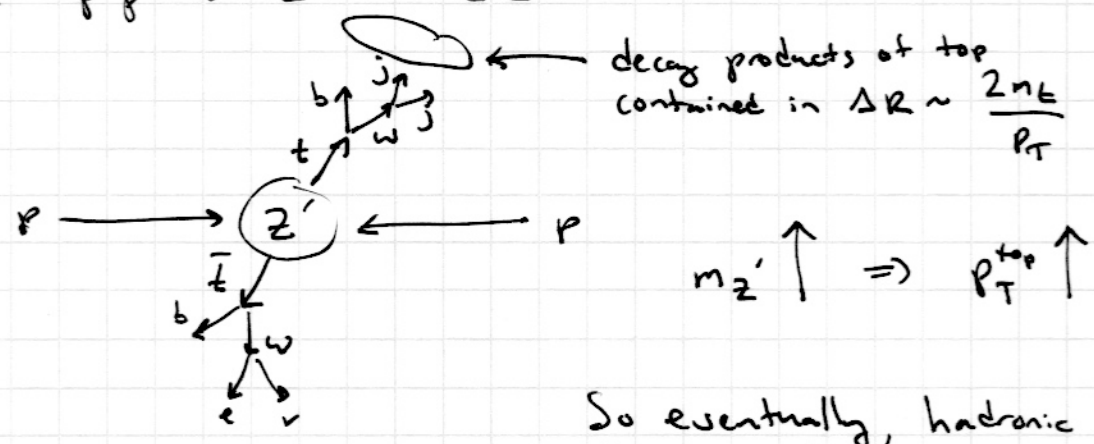
$$g \begin{array}{c} \diagup \\ \diagdown \\ \diagup \\ \diagdown \\ \diagup \\ \diagdown \end{array} \quad C_A = 3$$

Gluons have $9/4$ more radiation than quarks.

Putting aside questions about the intrinsic definition of "quark jet" and "gluon jet", I want to show you a calculation of how well you can discriminate these categories by measuring the jet mass. (Strictly speaking, a 2-point correlator.)



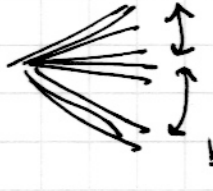

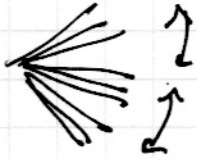
For additional context, jet classification became a hot topic once the new physics we were searching for became much heavier than $m_{W/Z/H/E}$.

E.g. $pp \rightarrow Z' \rightarrow t\bar{t}$



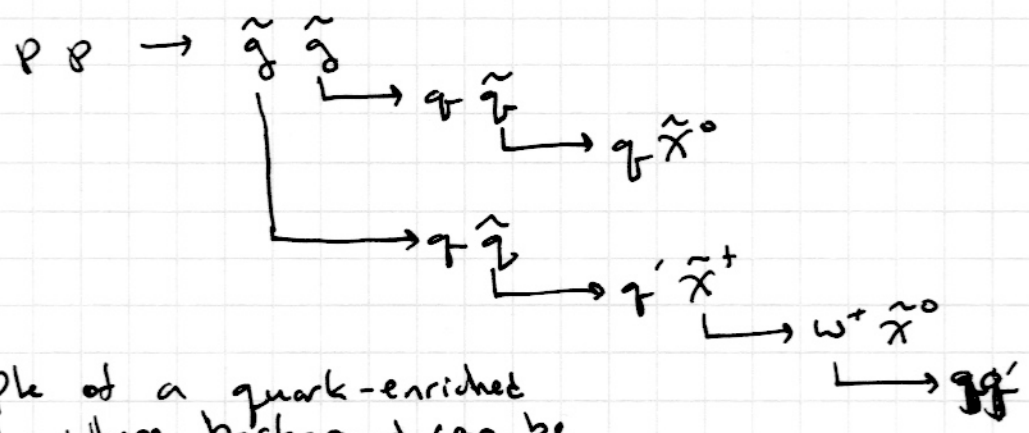
So eventually, hadronic top decay products were reconstructed as single 3-prong "fat jet"

By now, boosted jet tagging is standard.

- W/Z  2 prong
- Higgs  double b-tagged fat jet
- top  3 prong.
b-tag
- light Z'  2 prong again, different mass
- RPV gluino/neutralino  3 prong again, different mass.

In this context, quark/gluon discrimination is particularly challenging, since no "topological" distinction. Nevertheless, it is highly relevant for BSM searches.

E.g. gluino cascade decays



Example of a quark-enriched signal, where background can be gluon dominated.

We will focus on soft & collinear limit of QCD where FSR splitting probability is

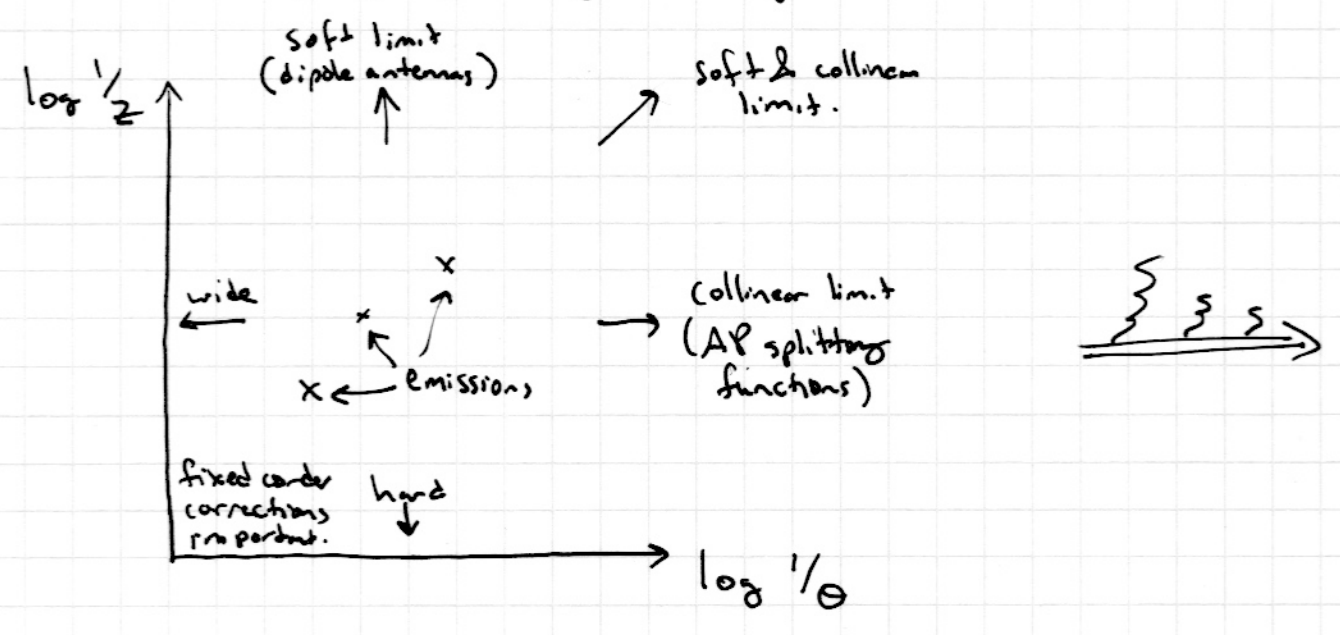
$$dP_{i \rightarrow ij} = \frac{2\alpha_s}{\pi} C_i \frac{dz}{z} \frac{d\Omega}{\Omega}$$

In this "strongly-ordered" (or "leading log" or "double log") limit, we can do simple but insightful calculations.

Very low accuracy, but gives a flavor for the types of questions you can ask/answer.

See SCET, direct resummation, ... for ways to systematically improve this picture.

Key insight: splitting kernel yields uniform emissions
in $(\log 1/\theta, \log 1/z)$ plane (sometimes called Lund plane)



Recursive application of $dP_{i \rightarrow jk}$ is how parton shower algorithms work, essentially ~~is~~ taking factorization to logical extreme.

This picture captures some information at all orders in d_s .
Often more realistic qualitatively than fixed-order methods.

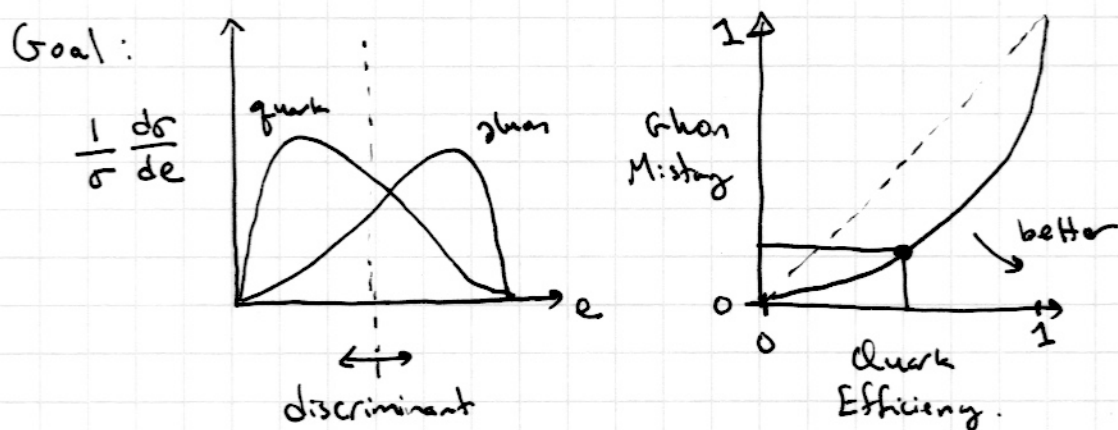
Main Insight: Jet is not just a single parton
(closer to Wilson-line-wrapped eikonal parton)

Very instructive to always keep in mind the soft gluon haze.

Concrete question: How well can you tag quark jets and reject gluon jets?

$C_F < C_A$, so gluon jets should be "fatter"

Need to define observables sensitive to this difference.



Ideally we would:

- Predict from first-principles QCD
- Validate in Monte Carlo parton showers
- Test in LHC data

Choice of discriminant observable?

- IRC safe (in order to perform perturbative calculations)
- Relatively insensitive to hadronization and other nonperturbative effects

$$\sigma_e \sim \left(\frac{\Lambda_{QCD}}{p_{T, \text{jet}}} \right)^{\#}$$

Here is a simple one-parameter family of quark/gluon discriminants
Energy-Energy Correlation Function

$$e_2^{(\beta)} = \frac{\sum_{ij} P_{T_i} P_{T_j} R_{ij}^\beta}{\left(\sum_i P_{T_i}\right)^2 R^\beta}$$

Angular exponent β should satisfy $\beta > 0$ for IRC safety
(Also known as $C_1^{(\beta)}$)
Convenient to normalize by jet radius

2-point correlator for 1-prong testing

Quick check for IRC safety:

- Soft safe? Yes, e_2 unaffected by $P_{T_i} \rightarrow 0$
- Collinear safe? Yes, e_2 is additive so $P_T \rightarrow P_{T,1} + P_{T,2}$ has no effect (for $\beta > 0$)

What is quark/gluon discrimination power for e_2 ?

We will work in extreme soft & collinear & strongly-ordered limit of QCD. (This is borderline physical, but good for intuition.)

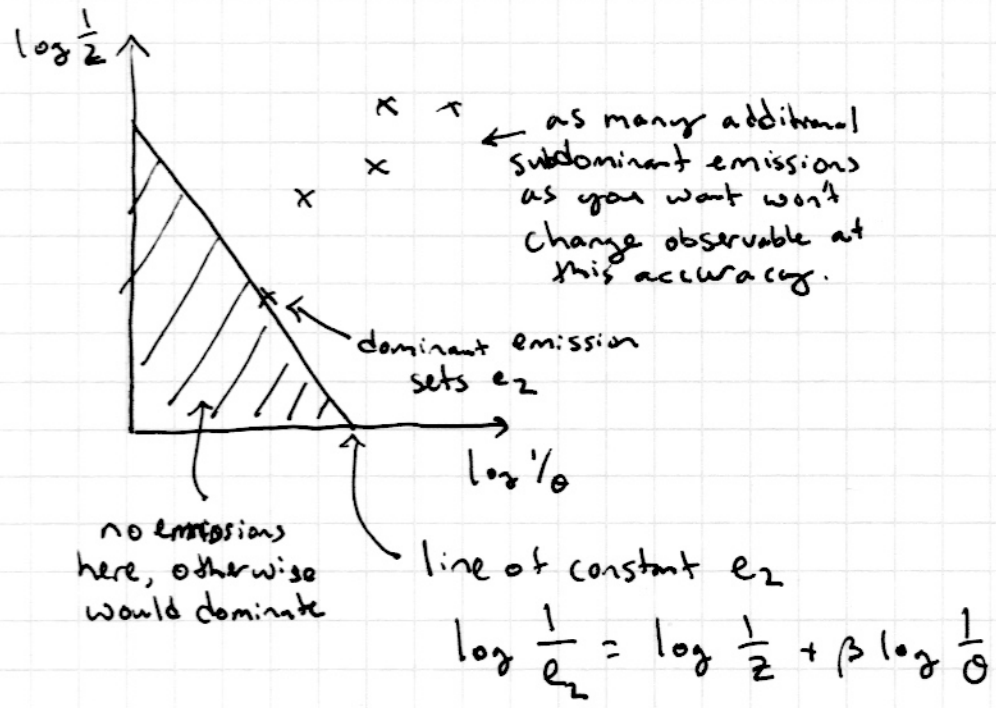


$$e_2 \approx \frac{P_{T1} P_{T2} R_{12}^\beta}{(P_{T1} + P_{T2})^2 R^\beta}$$

$$\approx z \Theta^\beta \quad \text{for } P_{T2} \ll P_{T1}$$

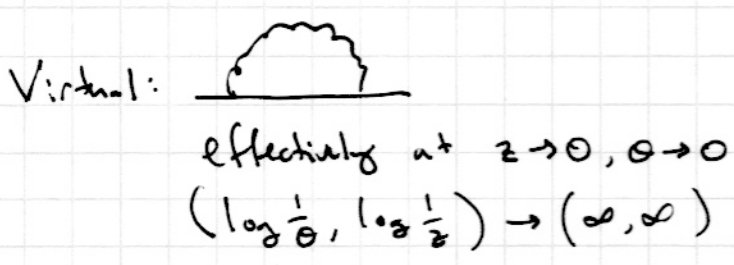
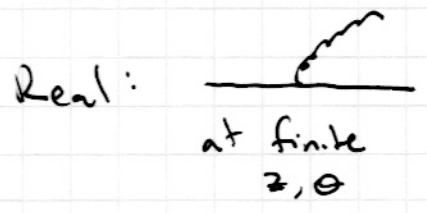
with $z \equiv \frac{P_{T2}}{P_{T1} + P_{T2}}$ $\Theta \equiv \frac{R_{12}}{R}$

We can plot this observable on the $(\log \frac{1}{\theta}, \log \frac{1}{z})$ plane.



Recall: $dP_{i \rightarrow j} = \frac{2\alpha_s}{\pi} C_i \frac{dz}{z} \frac{d\theta}{\theta}$, so we just have to compute probability to get specific value of e_2 , right?

But wait, $dP_{i \rightarrow j}$ is for real emissions, where are virtual diagrams?



Because our observable is IRC safe, we know divergences have to cancel, but how does this work?

Easiest way to think about this is in terms of probabilities

$$P_{\text{emit}} + P_{\text{no-emit}} = 1$$

At Born order:

$$\mathcal{O}(\alpha_s^0) \quad P_{\text{noemit}}^{(0)} = 1 \quad \xrightarrow{\text{just}}$$

At next order:

$$\mathcal{O}(\alpha_s^1) \quad P_{\text{emit}}^{(1)} + P_{\text{noemit}}^{(1)} = 0 \quad \uparrow \text{weird...}$$



Now we just have an exercise in probability!

Chance to get any value of e_2 less than e_2^{max} is

$$\sum_q \uparrow \text{for quark} \quad P(e_2^{\text{max}}) = 1 + \text{[trapezoid]} + \frac{1}{2} \text{[trapezoid]}^2 + \dots$$

\uparrow no emission at $\mathcal{O}(\alpha_s^0)$
 \uparrow all emissions below e_2^{max} at $\mathcal{O}(\alpha_s^1)$
 \uparrow symmetry factor
 \uparrow both emissions below e_2^{max} at $\mathcal{O}(\alpha_s^2)$

Then, using fact that $\text{[trapezoid]} = \text{[triangle]} + \text{[triangle]}$, we have:

\uparrow real plus virtual
 \uparrow just real

$$\sum_q \uparrow \text{for quark} \quad P(e_2^{\text{max}}) = \int_0^{e_2^{\text{max}}} p(e_2) de_2 = \exp \left[-\frac{2\alpha_s}{\pi} C_F (\text{area of } \text{[triangle]}) \right]$$

Probably most of you have never done a field theory calculation like that before!

Final result: $\Sigma_q(e_2^{\max}) = \exp \left[\frac{-\alpha_s}{\pi} \frac{C_F}{\beta} \log^2 \frac{1}{e_2^{\max}} \right]$

↑
called Sudakov form factor

↑
double logarithmic because of soft & collinear singularities.

For gluons, just swap $C_F \rightarrow C_A$.

This is a (baby) example of a resummed calculation where you capture some information to all orders in α_s . (can be systematically improved using approaches like SCET.)

To go from cumulative distribution Σ to probability distribution, just take a derivative.

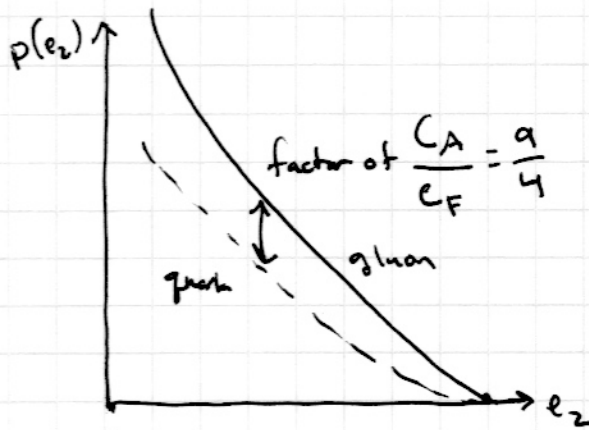
$$p(e_2) = \frac{1}{\sigma} \frac{d\sigma}{de_2} = \frac{\partial}{\partial e_2} \Sigma(e_2)$$

$$p_q(e_2) = \frac{2\alpha_s}{\pi} \frac{C_F}{\beta} \frac{1}{e_2} \log \frac{1}{e_2} \exp \left[-\frac{\alpha_s}{\pi} \frac{C_F}{\beta} \log^2 \frac{1}{e_2} \right]$$

at fixed order in α_s , this goes singular as $e_2 \rightarrow 0$

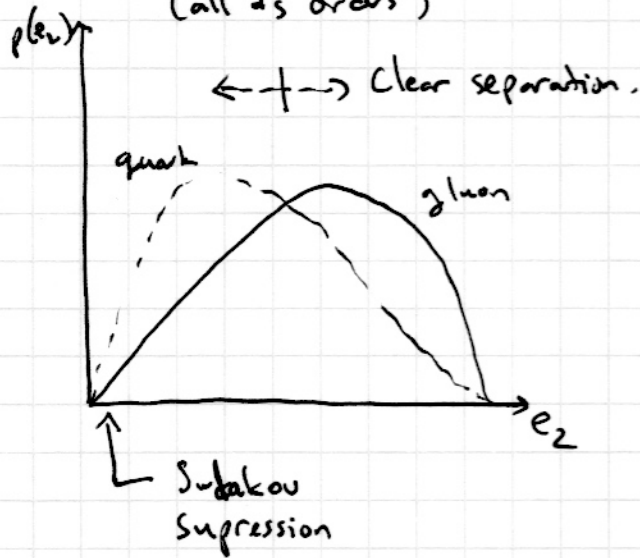
at all orders in α_s , you get Sudakov form factor that regulates singularity, giving much more physical result.

At fixed α_s order



plus/delta function from virtual diagrams to ensure normalized distributions

In strongly ordered limit (all α_s orders)



Finally, we can compute discrimination power.

Place a cut at e_2^{cut}

$$e_2 < e_2^{\text{cut}} \Rightarrow \text{"quark"}$$

$$e_2 > e_2^{\text{cut}} \Rightarrow \text{"gluon"}$$

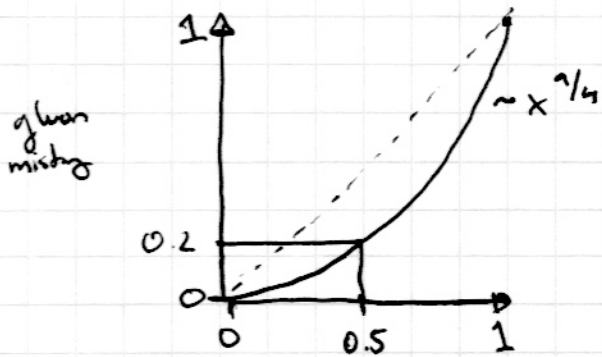
We want to know \uparrow mistag rate ($\sum_g(e_2^{\text{cut}})$)
as a function of quark efficiency ($\sum_q(e_2^{\text{cut}})$).

At this order, we have simple ~~the~~ "Casimir scaling"

$$\sum_g = \left(\sum_q \right)^{C_A/C_F} \rightarrow 9/4$$

(This simple relation is violated at higher orders. Note that at this order, it is independent of β .)

Final * mistag vs. efficiency curve:



Huge range of observables exhibit Casimir scaling, which is why quark/gluon discrimination is challenging.

But there are methods to achieve improved performance, which you can ask me about offline.

This calculation neglected many important higher-order effects.

- Multiple Emissions
- Color Coherence
- Subleading terms in AP splitting function
- Fixed-order corrections
- Running α_s
- Non-global logarithms
- Hadronization Effects
- Underlying event contamination
- ...

Mostly, I hope you've gained some appreciation of why we need to define observables and why it is beneficial to work to all orders in α_s .

Concluding Thoughts

Here are three things I want you to remember about QCD and collider physics.

- ① Have to think about observables. Just knowing scattering amplitudes is not enough to make predictions.

Right now, there is no "theory of observables", so you ~~are~~ have to assess things case by case.

- ② Factorization is crucial for making predictions.

You can't predict where every pion goes. Need to choose observables that respect, e.g., PDFs, otherwise you can't (yet) make first-principles predictions.

- ③ You see (quasi-)stable particles in your detector, not Standard Model ~~states~~ states.

Good observables on hadrons yield good proxies for QCD partons. Jet algorithms are one way to construct such proxies.

If you want to use perturbative proxies, need to use observables that are infrared and collinear safe.

Finally, I started these lectures with the master formula: (60)

$$\sigma_{\text{obs}} = \frac{1}{2E_{\text{cm}}^2} \sum_{n=2}^{\infty} \int d\mathbb{I}_n |M_{AB \rightarrow 12 \dots n}|^2 f_{\text{obs}}(\mathbb{I}_n)$$

Are there things you can do with collider data that don't fall into this framework?

Yes! There are data analysis techniques (e.g. from machine learning) that aren't based on histograms.

In many ways, frontiers of collider physics will involve thinking about the whole "space of measurements" and striking right balance between experimental robustness and theoretical calculability.

I hope I've given you some sense of the richness of QCD and collider physics in these lectures!