# Jet Templates

## Searches for High Multiplicity New Physics

## Jay Wacker

Quora

Boost 2014

August 18, 2014

w/ E. Izaguirre
A. Hook
M. Lisanti
S. El Hedri
M. Jankowiak
T. Cohen
H. Lou

# New Physics Searches

Rely heavily on one object
that QCD doesn't directly produce

Gives parametric control of QCD background

# Why are we waiting for discovery?

Signals could be just out of reach

Is there something that we're missing?

One dark corner:
Hadronic Final States

Missing usual
handles to control
QCD
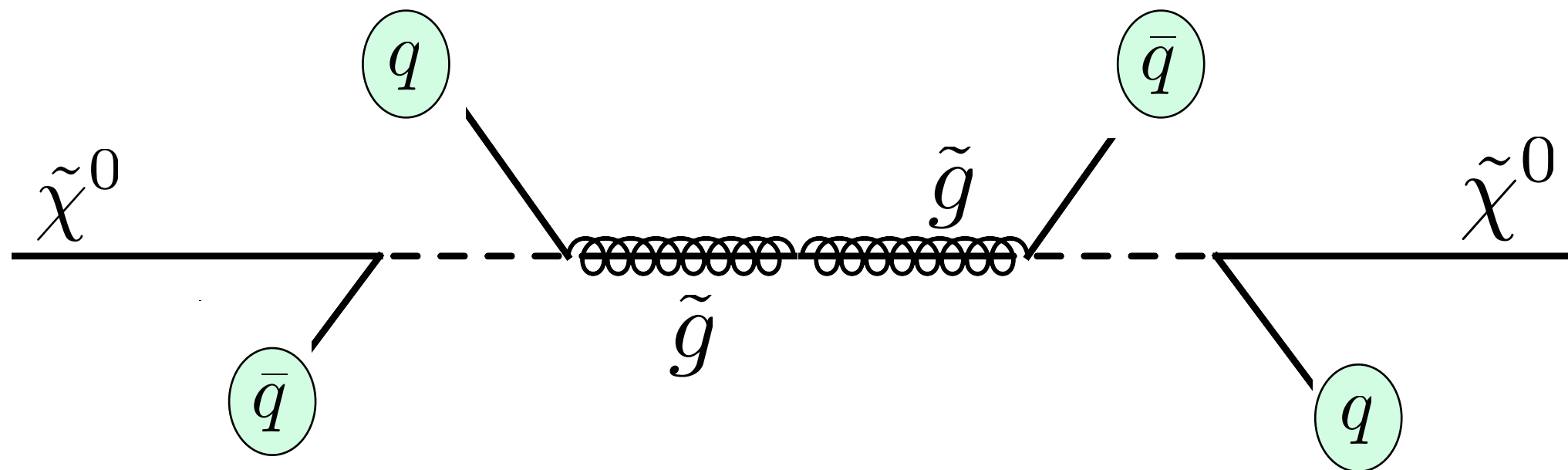
# Baryonic R-Parity Violation

Eviscerates MET

$$\int d^2\theta \;\; \lambda''_{ijk} \, U_i^c D_j^c D_k^c$$

## Makes LSP decay

to 3 quarks (most LSPs)
to 2 quarks (squark LSPs)
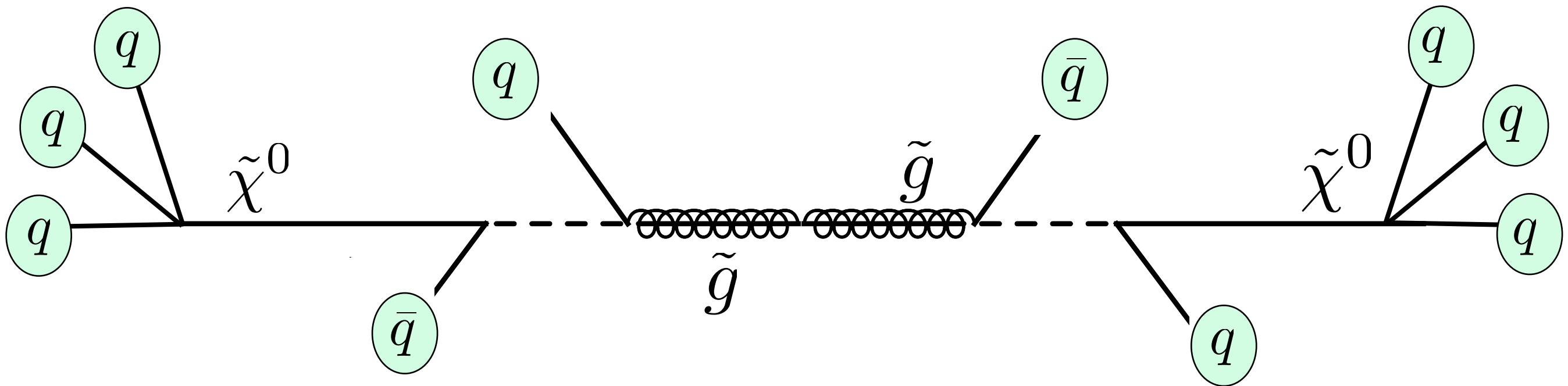
(one quark could be top → +2j)

# Increases multiplicity significantly
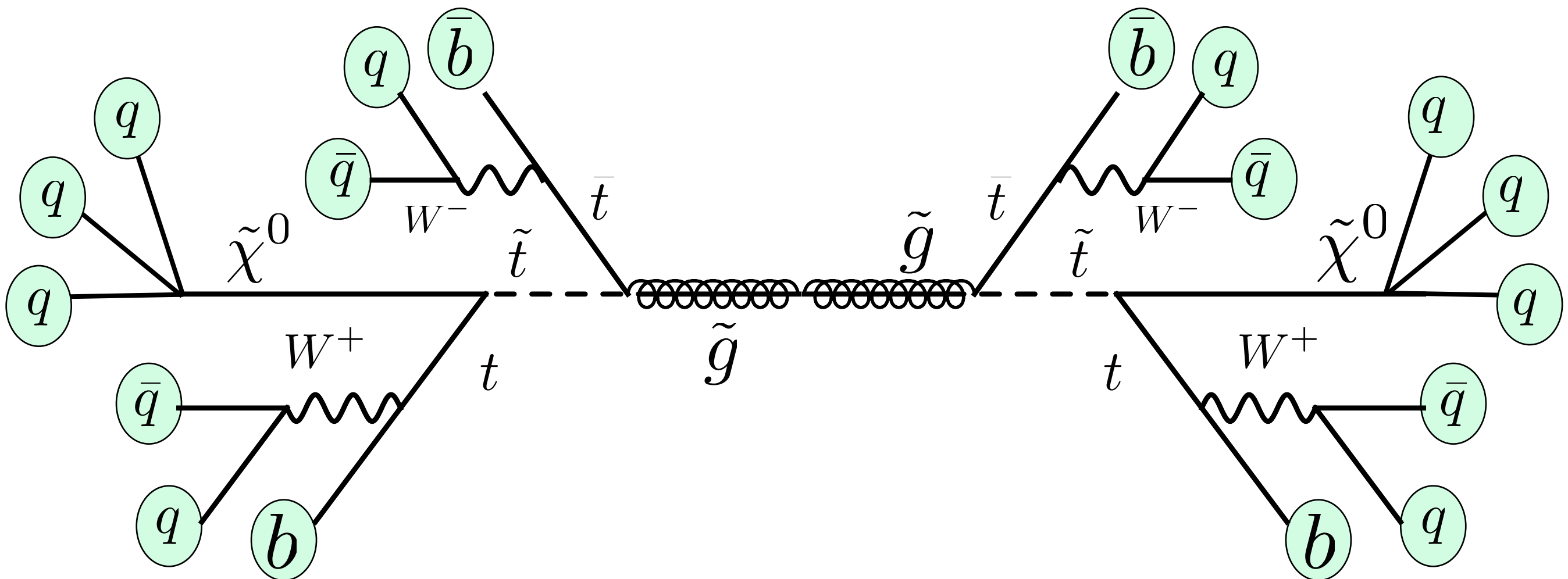
# The          Classic Susy Signature

# The Less-Classic Susy Signature
# 10⁺ Partons no MET

# The Less-Classic Natural Susy Signature

# 18⁺ Partons



Still some MET from W decays, but much less
Don't want to pay SSDL branching ratio (lepton isolation is hard)

# Main Point:

Many signals of new physics
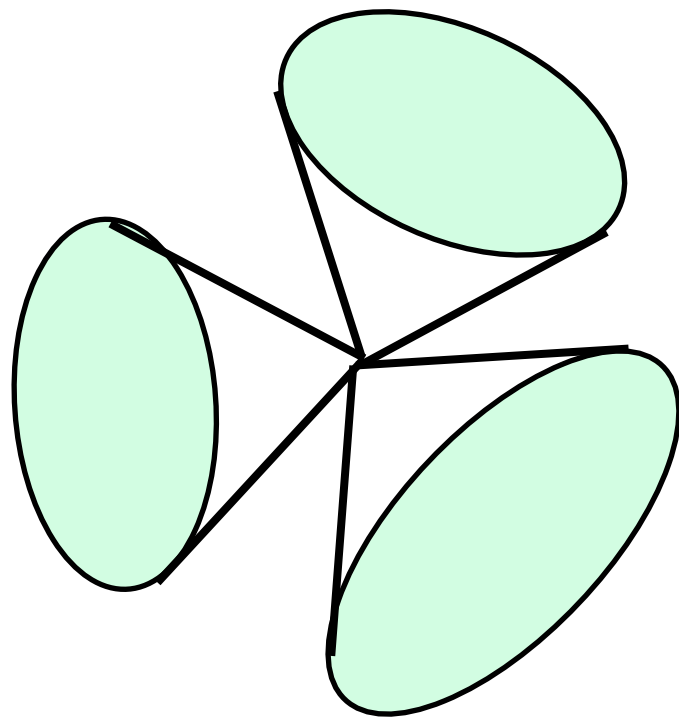produce lots of final state quarks or gluons

Easy to come up with other signals
with high multiplicity signals

Don't want to have a dedicated
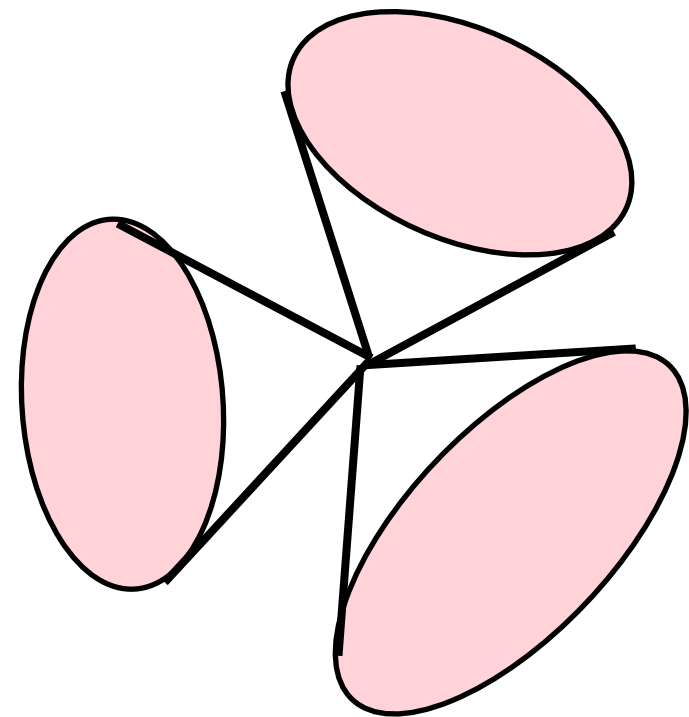search for every possibility

Want to use the multiplicity to distinguish
SM from BSM

# Need a handle to distinguish

Normal QCD Multijet

BSM Multijet

# Fat Jets

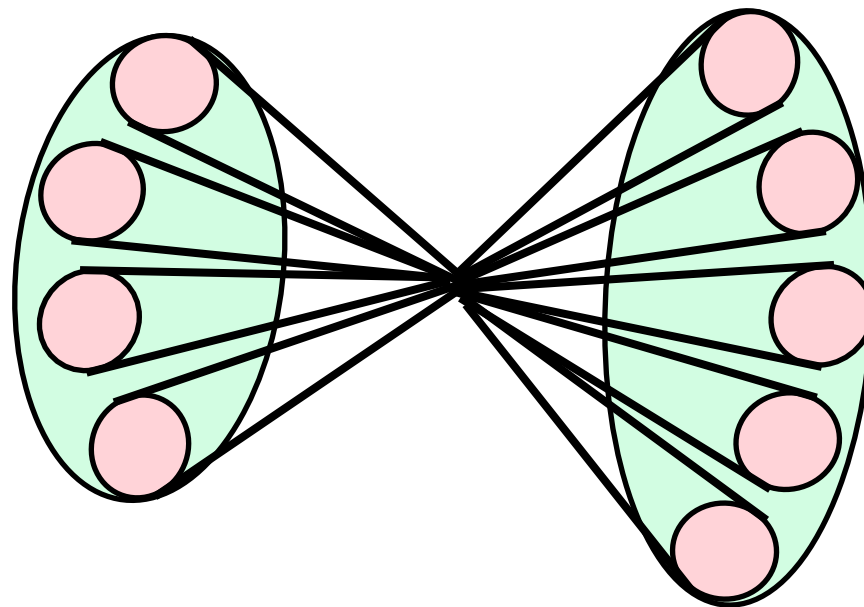## Fat Jets Coarse Grain the Phase Space

Easy to construct inclusive kinematic signals using fat jets

Thin Jets are great at determining multiplicity,
but constructing meaningful variables
out of a heterogeneous high dimensionful space is hard

Identify high multiplicity based
upon Fat Jet observables

# Truth Of QCD Multijets
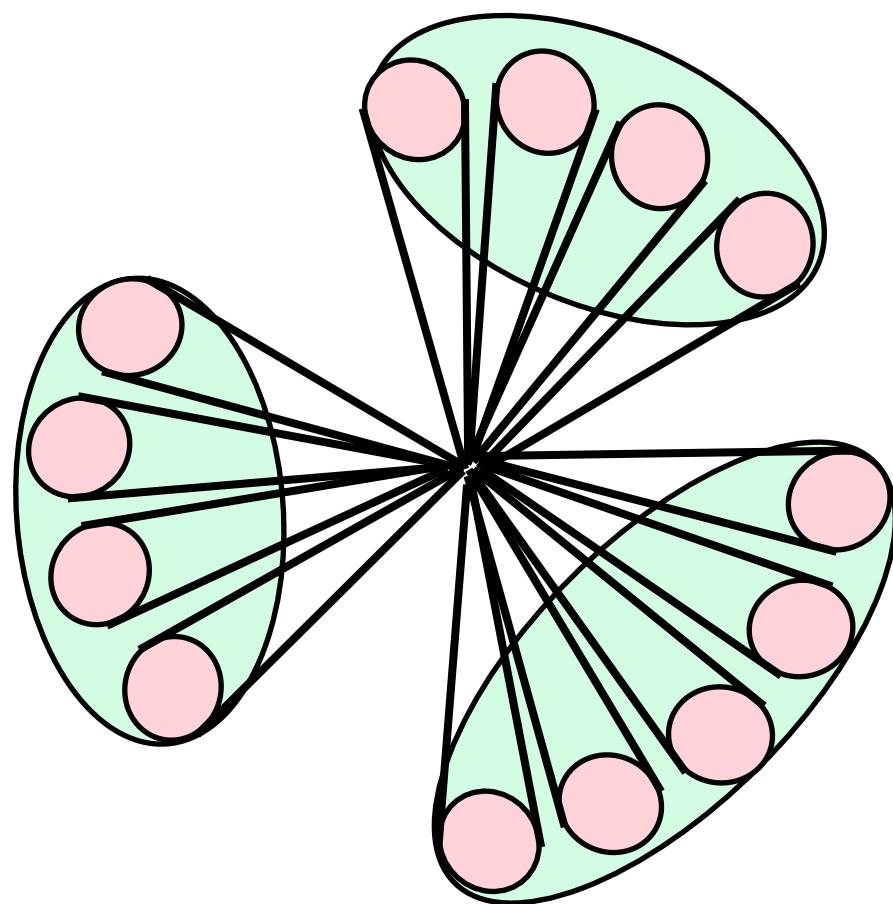
## Many QCD Multijets are glorified Dijets



Requiring 3 or 4 Fat Jets is a serious reduction
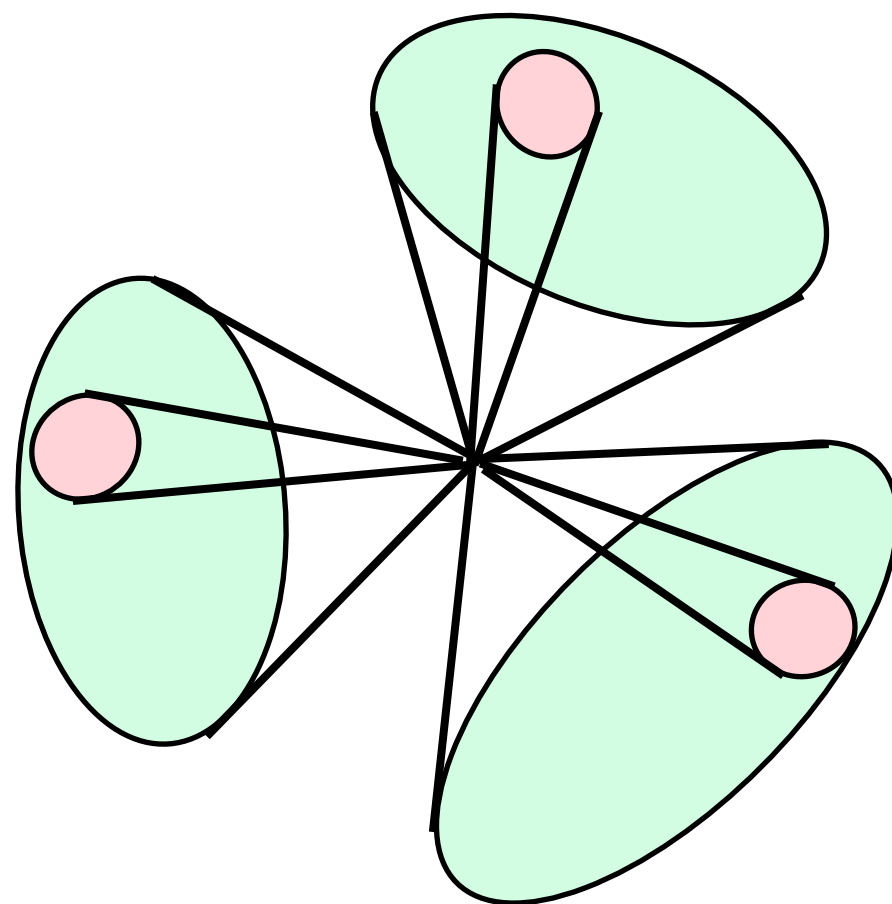in QCD rate

4 Fat jets is really a 2 → 4 process
6 Thin jets is dominated by 2 → 2 + parton showering

# Still need to distinguish

Signal

Background

# The difference between them is clear

## Large Invariant Mass

## Small Invariant Mass

$$\frac{m_j}{p_T} \sim 1$$

$$\frac{m_j}{p_T} \sim 0.3$$

More jet substructure

Less jet substructure

# Introduce Jet Observables

Sum of Jet Masses

$$M_J = \sum_{n=1}^{N_J} m_{j_n}$$

QCD jets have most of their mass generated by the parton shower

Top events have their mass capped near 400 GeV

# Subjettiness

Jet mass is the coarsest measure of jet substructure

Equal pT and mass jets

 versus 

Massive QCD jets mostly have 2 subjets

High multiplicity signals are more subjets

Used kT method of counting subjets

(1302.1870)

More than a Mass Cut

Fraction of Jets with $N_{subjets}$

$p_T = 100 - 200$ GeV

Legend: 1, 2, 3, 4

$m/p_T$

4 Fat Jets, $p_T > 100$ GeV

After $\not{E}_T > 150$ GeV

$M_J$ Distribution

Top
V+jets
Diboson
Misc
QCD

$\sigma(fb)/25$GeV

$M_J$ (GeV)

4 Fat Jets, $p_T > 100$ GeV

After $\not{E}_T > 150$ GeV & $M_J > 280$ GeV
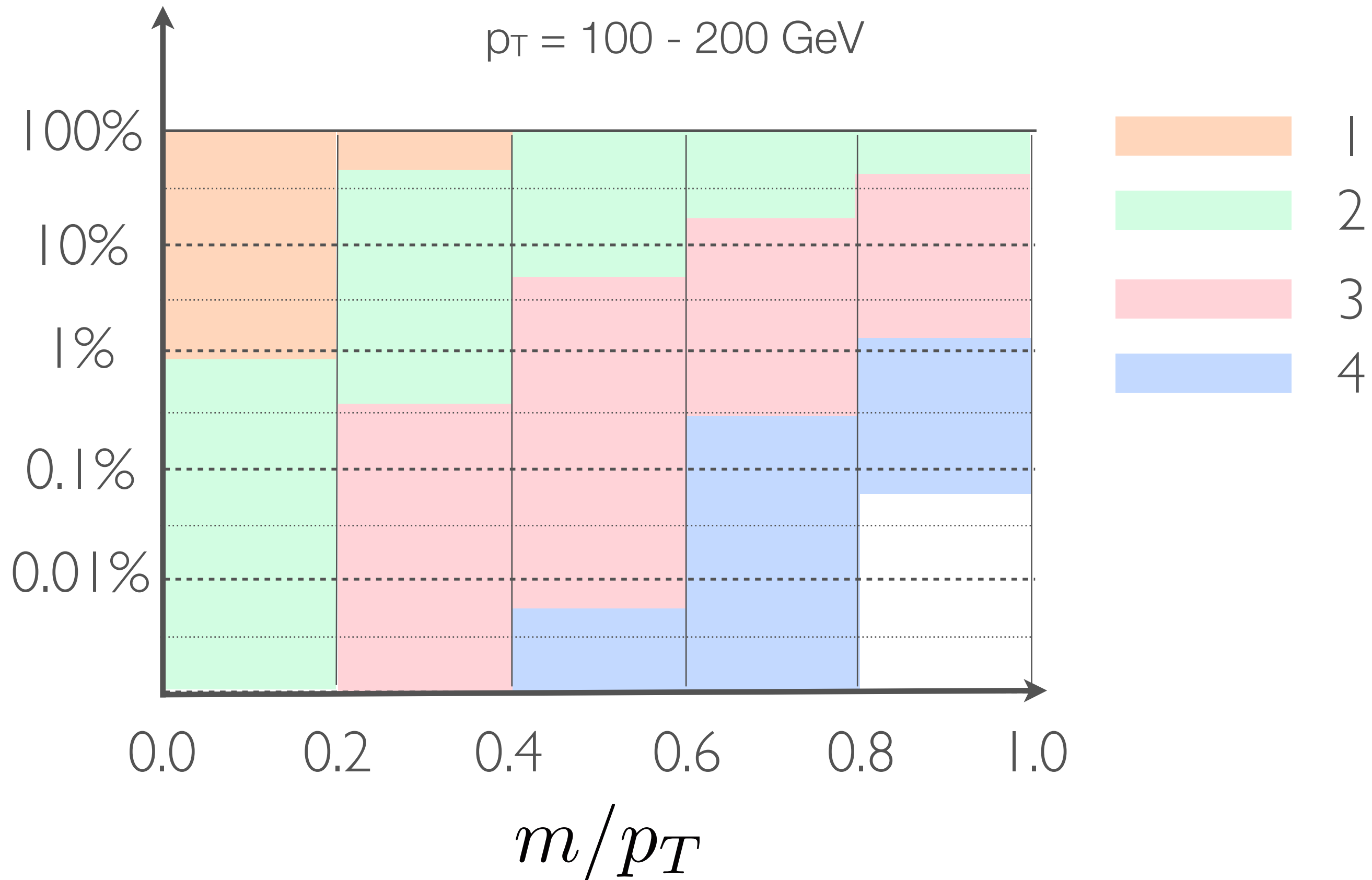
N$_J$ Distribution

# Improvements of $N_J$ vs $M_J$ only Search

$$E\!\!\!/_T > 125 \text{ GeV} \qquad M_J \geq 425 \text{ GeV} \qquad N_J > 14$$

A little bit of MET from W-decays



$\mathcal{G}_7$

30fb$^{-1}$ 8 TeV

| | |
|---|---|
| - - - | $M_J + E\!\!\!/_T$ |
| —— | $M_J + E\!\!\!/_T + N_{CA}$ |
| -·-·- | ATLAS |
| - - - | CMS/5 |

$\sigma(\tilde{g}\,\tilde{g})_{NLO}$

$\sigma \times Br$ (fb)

$m_{\tilde{g}}$( GeV)

$\sigma_{SM} \simeq 0.07\text{fb}$

Factor of 8 improvement in cross section, factor of 64 less luminosity

# Variables are Great

… but Monte Carlos can't reproduce
all of jet substructure

## How to get backgrounds?

Particularly challenging when
variables are correlated

# Jet Factorization

QCD jets only have small correlations

Data driven background predictions possible



$$x = m_j/p_T$$

$$P_3(x_1, x_2, x_3) \simeq P_1(x_1)P_1(x_2)P_1(x_3)$$

$P_1$: Probability of a jet with m/$p_T$ = x

$P_3$: Probability of getting 3 jets with $x_1$, $x_2$, $x_3$

Measure in one sample and extrapolate

Also can use other control regions (MET/leptons/bjets)

# Natural "Data-Driven" approach to backgrounds

Measure $P_1(x; p_T)$ in dijets, use in multijets

Predict event-by-event acceptances

(probability an event passes cut)

$$A(p_{T1}, p_{T2}, p_{T3}) = \int_{M_J > m_{\mathrm{cut}}} d^3x \quad P_1(x_1; p_{T1}) P_1(x_2; p_{T2}) P_1(x_3; p_{T3})$$

Differential acceptance rate as a function of the kinematic variables

Can make an $M_J$ prediction based upon the events *measured*

Don't need to be able to calculate $M_J$ distribution
from first principles

# The Basic Idea of Jet Templates

# More Formally

k are kinematic variables          x are substructure variables

$$\frac{\mathrm{d}^{2N_j}\sigma(\vec{x}_i, \vec{\mathrm{k}}_i,)}{\mathrm{d}\vec{x}_1...\mathrm{d}\vec{x}_{N_j}\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}}$$

# More Formally

k are kinematic variables

x are substructure variables

$$\frac{\mathrm{d}^{2N_j}\sigma(\vec{x}_i, \vec{\mathrm{k}}_i,)}{\mathrm{d}\vec{x}_1...\mathrm{d}\vec{x}_{N_j}\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}} = \frac{\mathrm{d}^{N_j}\sigma(\vec{\mathrm{k}}_i)}{\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}} \, \rho(\vec{x}_1, ..., \vec{x}_{N_j} | \vec{\mathrm{k}}_1, ..., \vec{\mathrm{k}}_{N_j})$$

# More Formally

k are kinematic variables

x are substructure variables

$$\frac{\mathrm{d}^{2N_j}\sigma(\vec{x}_i, \vec{\mathrm{k}}_i,)}{\mathrm{d}\vec{x}_1...\mathrm{d}\vec{x}_{N_j}\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}} = \frac{\mathrm{d}^{N_j}\sigma(\vec{\mathrm{k}}_i)}{\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}} \rho(\vec{x}_1,...,\vec{x}_{N_j}|\vec{\mathrm{k}}_1,...,\vec{\mathrm{k}}_{N_j})$$

Approximate the multivariate joint distribution function
as independent distribution functions

$$\frac{\mathrm{d}^{N_j}\sigma(\vec{\mathrm{k}}_i)}{\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}} \rho(\vec{x}_1,...,\vec{x}_{N_j}|\vec{\mathrm{k}}_1,...,\vec{\mathrm{k}}_{N_j}) = \frac{\mathrm{d}^{N_j}\sigma(\vec{\mathrm{k}}_i)}{\mathrm{d}\vec{\mathrm{k}}_1...\mathrm{d}\vec{\mathrm{k}}_{N_j}} \prod_{i=1}^{N_j} \rho_i(\vec{x}_i|\vec{\mathrm{k}}_i).$$

# MEASURING THE TEMPLATES

Getting the central value is easy
Getting error bars is hard

## Used Kernel Smoothing

Take every event and replace its properties
with a Gaussian

$$\rho(m) = \sum_i \delta(m - m_i) \rightarrow \sum_i \exp\left(-\frac{(m - m_i)^2}{\sigma^2}\right)$$
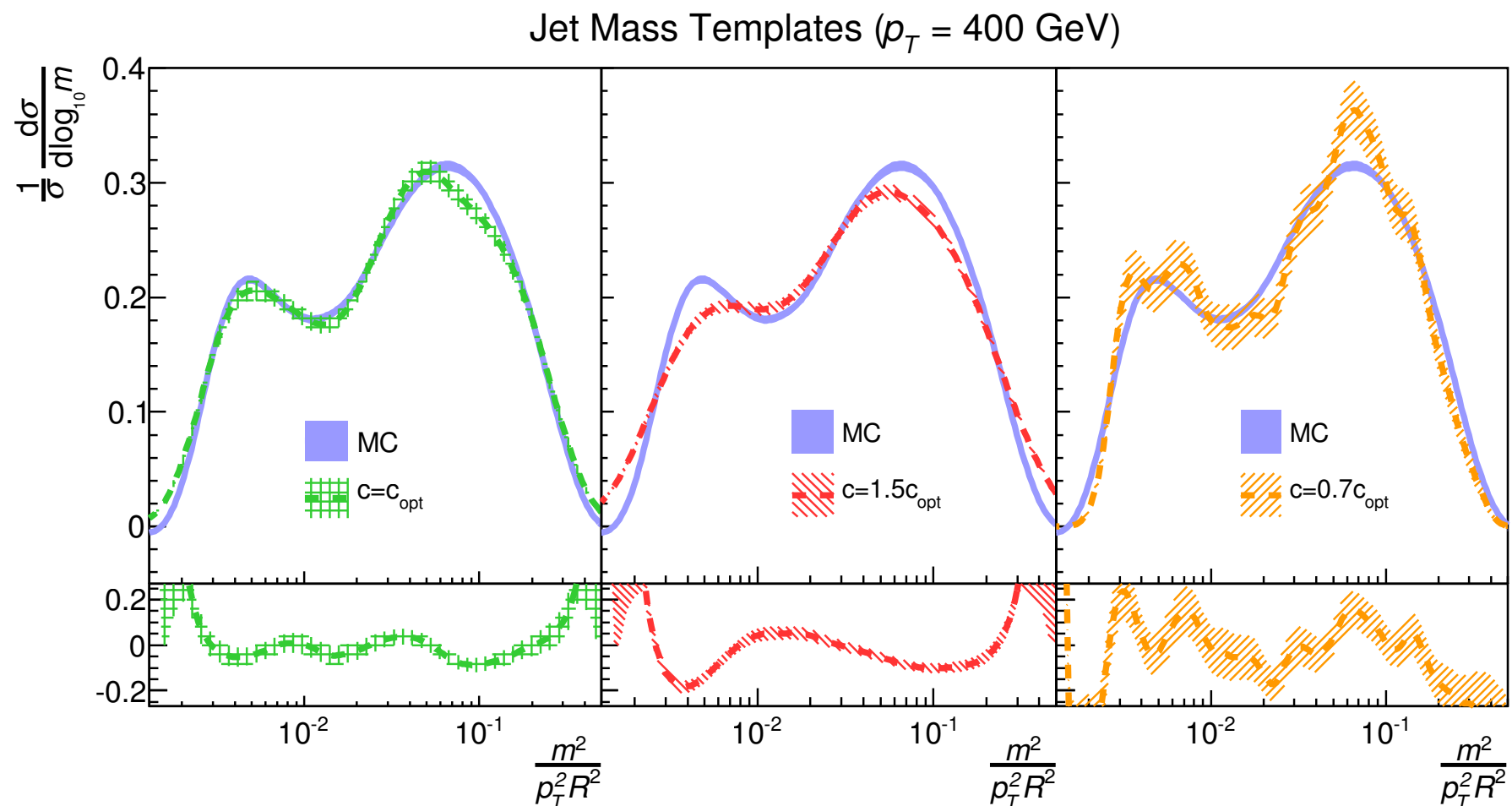
What is σ ?

# CHOOSING THE BANDWIDTH

Two separate errors arise in any procedure like this

## Variance & Bias

If you choose σ too small, then there
is a lot of statistical noise

If you choose σ too big, then there
the distribution systematically moves away from the true one



Jet Mass Templates ($p_T$ = 400 GeV)

# OPTIMAL BANDWIDTH

Typically chosen by "AMISE"
(asymptotic mean integrated square error)

$$\mathrm{AMISE}(\sigma) = \int dm \, \Big( \rho_0(m) - \rho(m; \sigma) \Big)^2$$

Can prove lots of things about this

# OPTIMAL BANDWIDTH

Typically chosen by "AMISE"
(asymptotic mean integrated square error)

$$\mathrm{AMISE}(\sigma) = \int dm \left( \rho_0(m) - \rho(m;\sigma) \right)^2$$

Can prove lots of things about this

But minimizing this is not the right thing to do

Variance is a Gaussian distribution

Bias is not, has non-Gaussian tails

# OPTIMAL BANDWIDTH

Typically chosen by "AMISE"

(asymptotic mean integrated square error)

$$\text{AMISE}(\sigma) = \int dm \, \Big( \rho_0(m) - \rho(m; \sigma) \Big)^2$$

Can prove lots of things about this

But minimizing this is not the right thing to do

Variance is a Gaussian distribution

Bias is not, has non-Gaussian tails
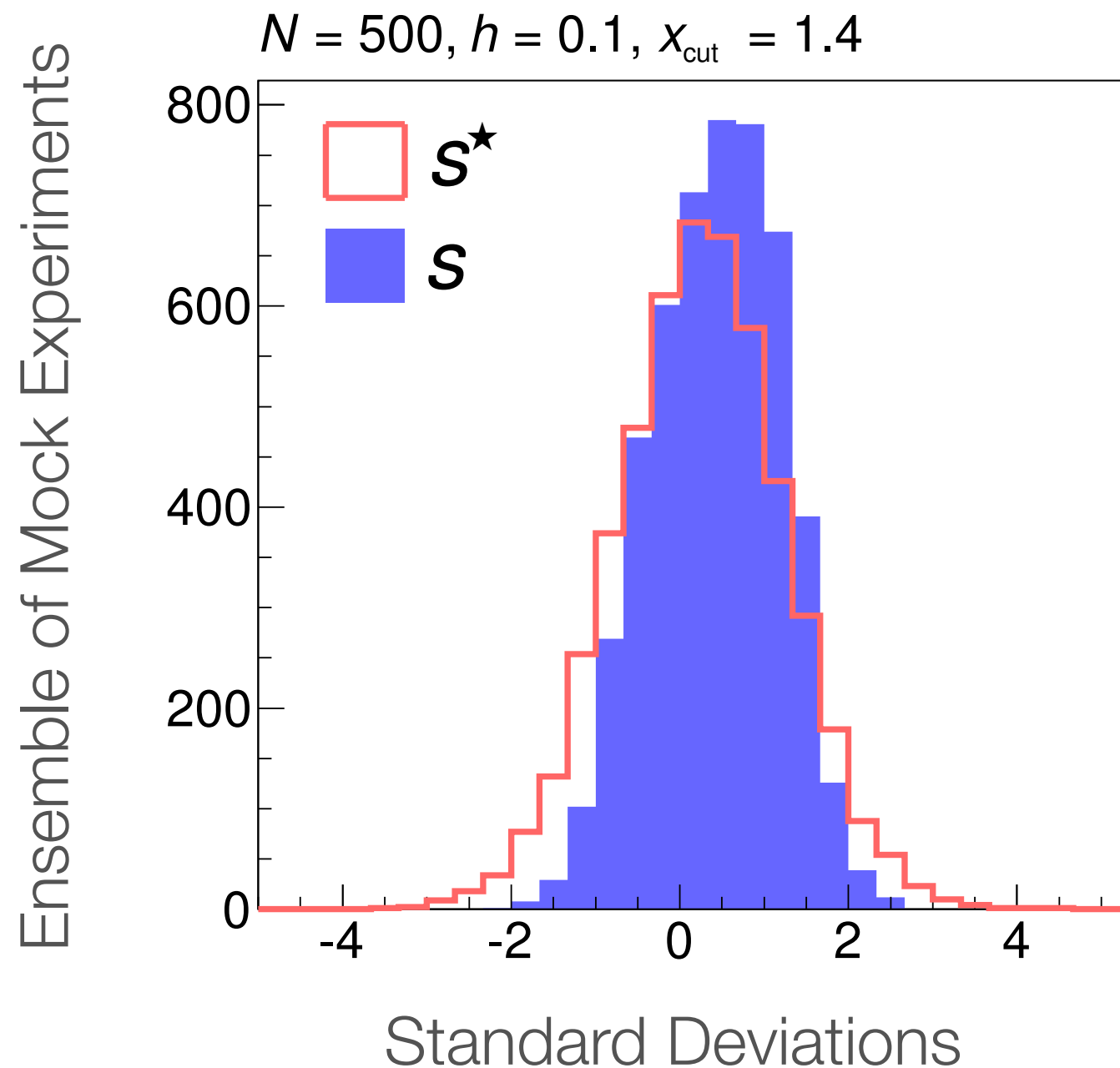
Want Variance to dominate over Bias

AMISE is a relatively function of bandwidth

Want to "undersmooth" the distribution

# BIAS-CORRECTED TEMPLATES
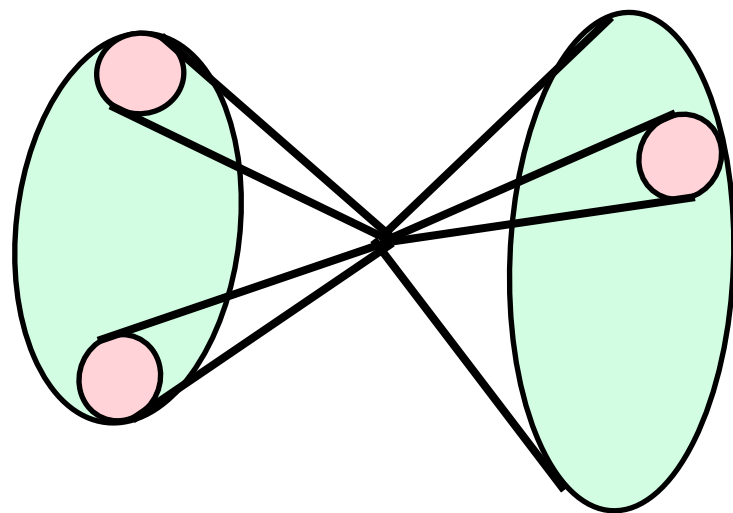
## Can measure the bias and correct for it at leading order

### Distributions are Gaussian, with width 1 and centered at 0
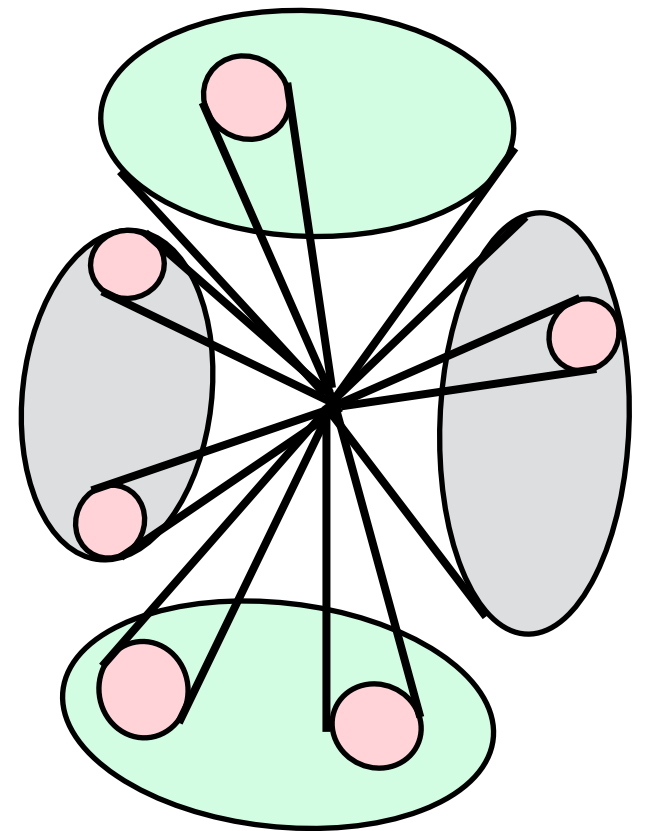
# Explicit Validation

## Control Region

Exclusive 2-Jets Events



## Signal Region

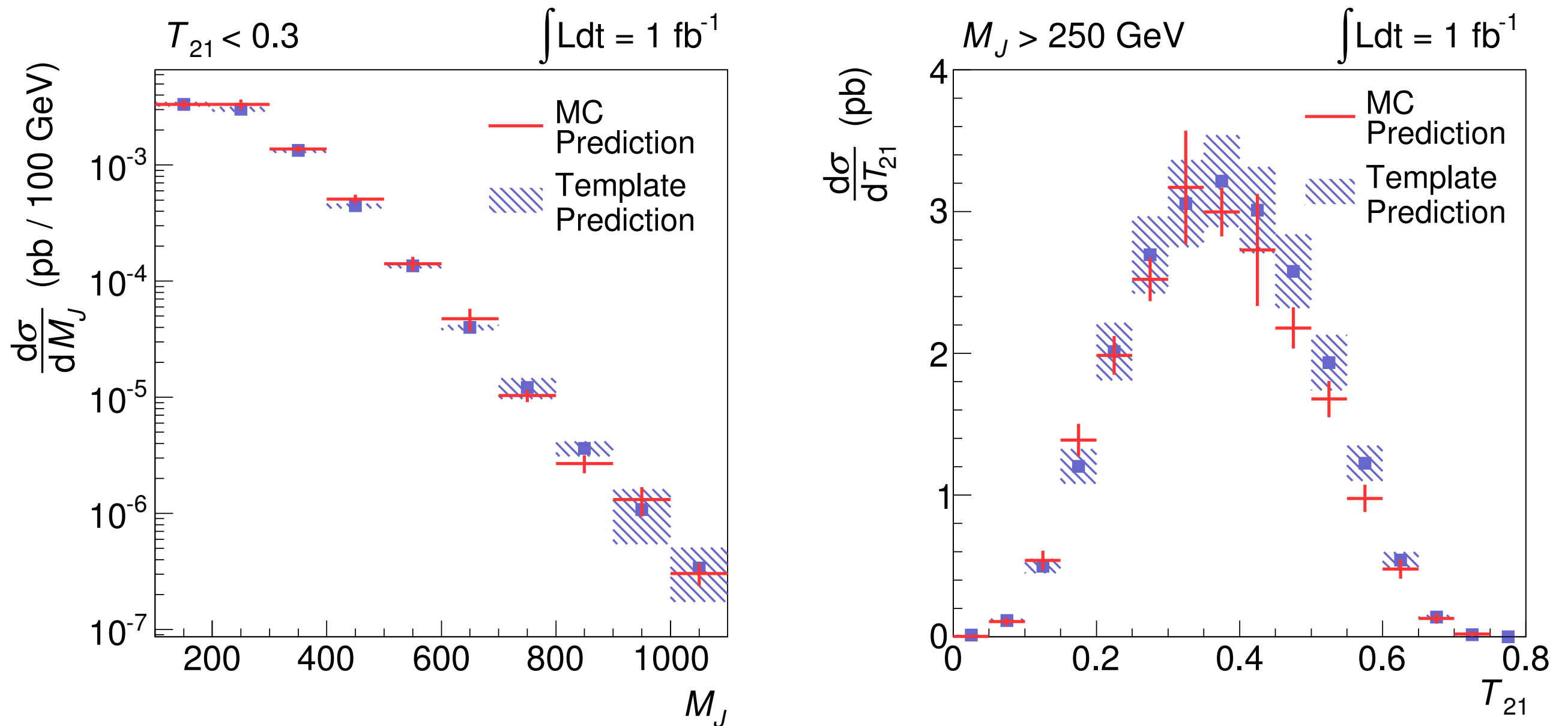Leading 2 Jets of 4-Jet Events



Test 2 Variables

$$M_J = m(j_1) + m(j_2)$$

$$T_{21}{}^2 = \tau_{21}(j_1) \ \tau_{21}(j_2)$$

# Works well in Monte Carlo

Take Exclusive Dijets and apply it to leading 2 jets in 4-Jet events

< 10% systematic differences



Minimally, jets in MC have less information,
can get more mileage with smaller MC calculations

# Works similarly well in Search Regions

$$\hat{\rho}^{\star} = \hat{\rho}^{\star}\left( -\log_{10}\left(\frac{m}{p_T}\right), \ \tau_{21}, \ \ln\left(\frac{p_T}{200 \text{ GeV}}\right)\right)$$

| $c$ | $M_J$ cut [GeV] | $T_{21}$ cut | MC | Template $\pm \hat{\sigma}_V \pm \hat{\sigma}_B$ |
|---|---|---|---|---|
| 0.37 | 500 | 0.3 | $20.3 \pm 2.2$ | $19.2 \pm 2.3 \pm 0.6$ |
| 0.52 | 750 | 0.3 | $0.86 \pm 0.10$ | $0.96 \pm 0.19 \pm 0.05$ |
| 0.37 | 500 | 0.6 | $45.8 \pm 3.5$ | $45.2 \pm 3.7 \pm 1.3$ |
| 0.52 | 750 | 0.6 | $1.67 \pm 0.14$ | $1.90 \pm 0.19 \pm 0.13$ |

Always under-smoothed to make the calculated bias smaller than the expected variance dominate
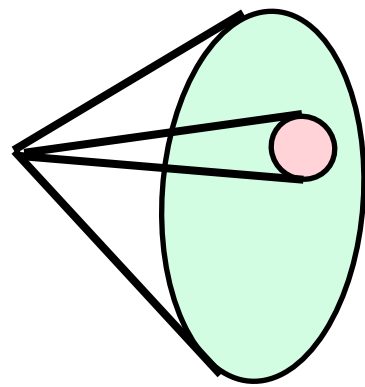
# Did this have to work?

No! A non-trivial check

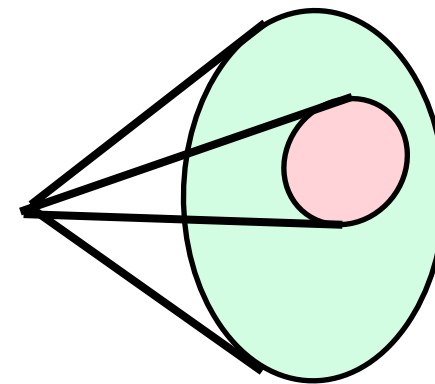For instance, Quark vs Gluon Jets

### Quarks:
Smaller Color, Less radiation

### Gluons:
Bigger Color, More radiation



Full Dijet Sample is

$$\rho_{12}(\vec{x}_1, \vec{x}_2) = c_{qq}\rho_{qq}(\vec{x}_1, \vec{x}_2) + c_{qg}\rho_{qg}(\vec{x}_1, \vec{x}_2) + c_{gq}\rho_{gq}(\vec{x}_1, \vec{x}_2) + c_{gg}\rho_{gg}(\vec{x}_1, \vec{x}_2),$$
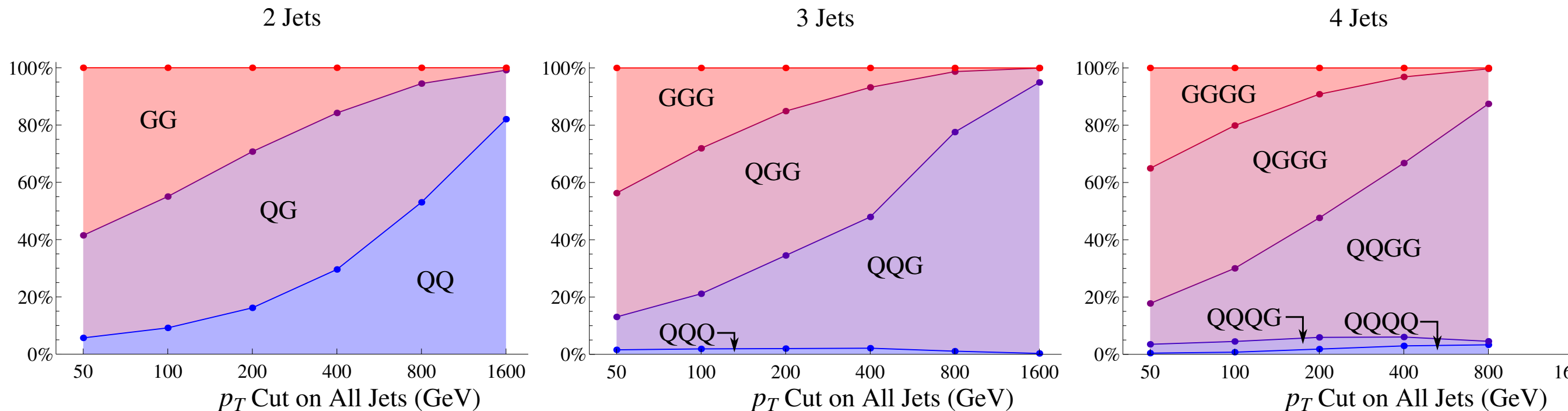
Approximating by

$$\tilde{\rho}(\vec{x}_1, \vec{x}_2) = \tilde{\rho}(\vec{x}_1)\tilde{\rho}(\vec{x}_2)$$

$$\tilde{\rho}(\vec{x}) = \left( c_{qq} + \frac{c_{qg} + c_{gq}}{2} \right) \rho_q(\vec{x}) + \left( c_{gg} + \frac{c_{qg} + c_{gq}}{2} \right) \rho_g(\vec{x}).$$

# Desperately Seeking Correlations

Have seen no evidence yet of correlations

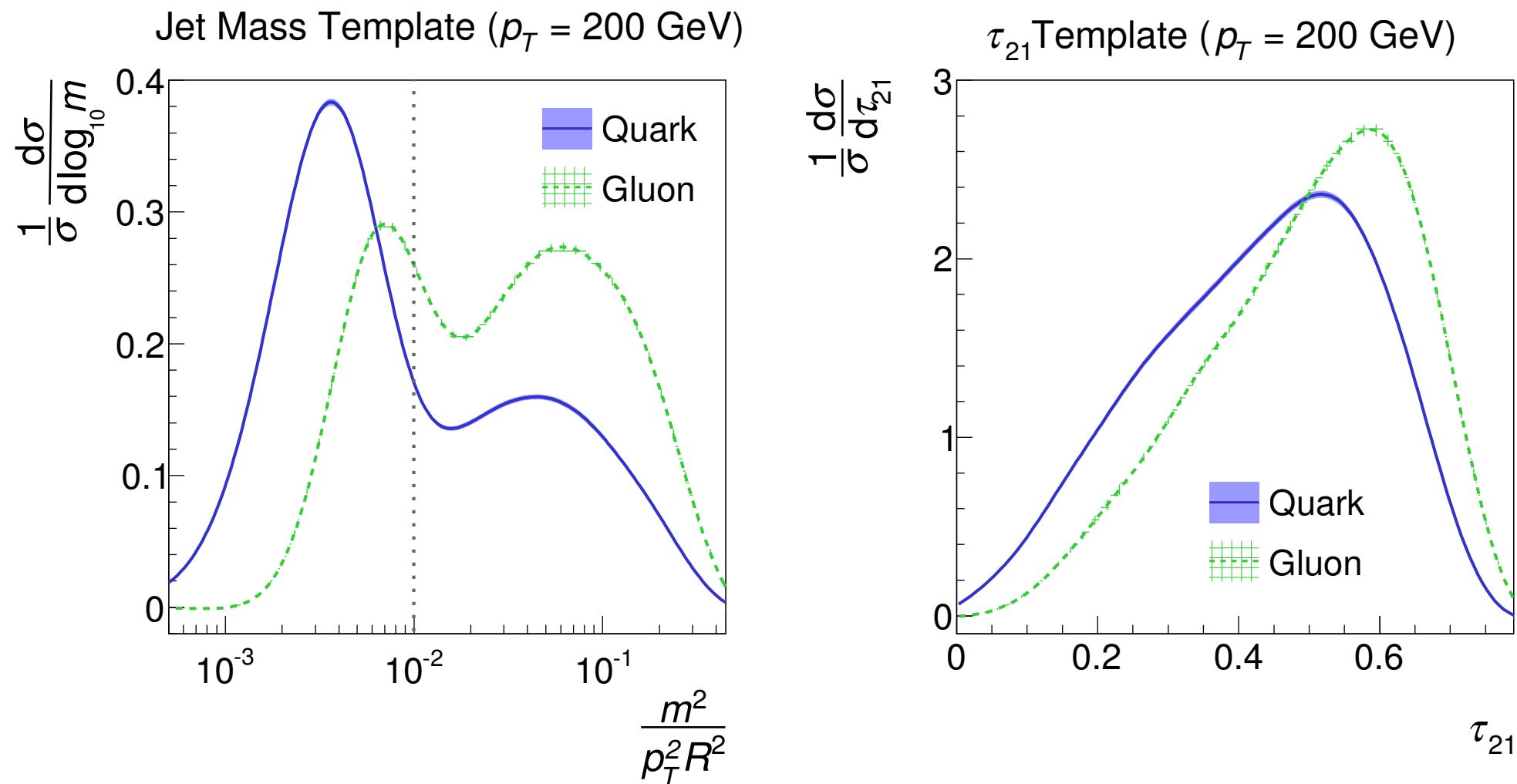Look at samples with different compositions



Leading 2 jets similar enough in composition between 2Jets & 4Jets

Using single template on all 4 jets doesn't work

# Q vs G Distributions Are Different

Have similar shapes and compositions cancel



Jet Mass Template ($p_T$ = 200 GeV)

$\frac{1}{\sigma}\frac{d\sigma}{d\log_{10} m}$

— Quark
— Gluon

$\frac{m^2}{p_T^2 R^2}$

$\tau_{21}$ Template ($p_T$ = 200 GeV)

$\frac{1}{\sigma}\frac{d\sigma}{d\tau_{21}}$

— Quark
— Gluon

$\tau_{21}$

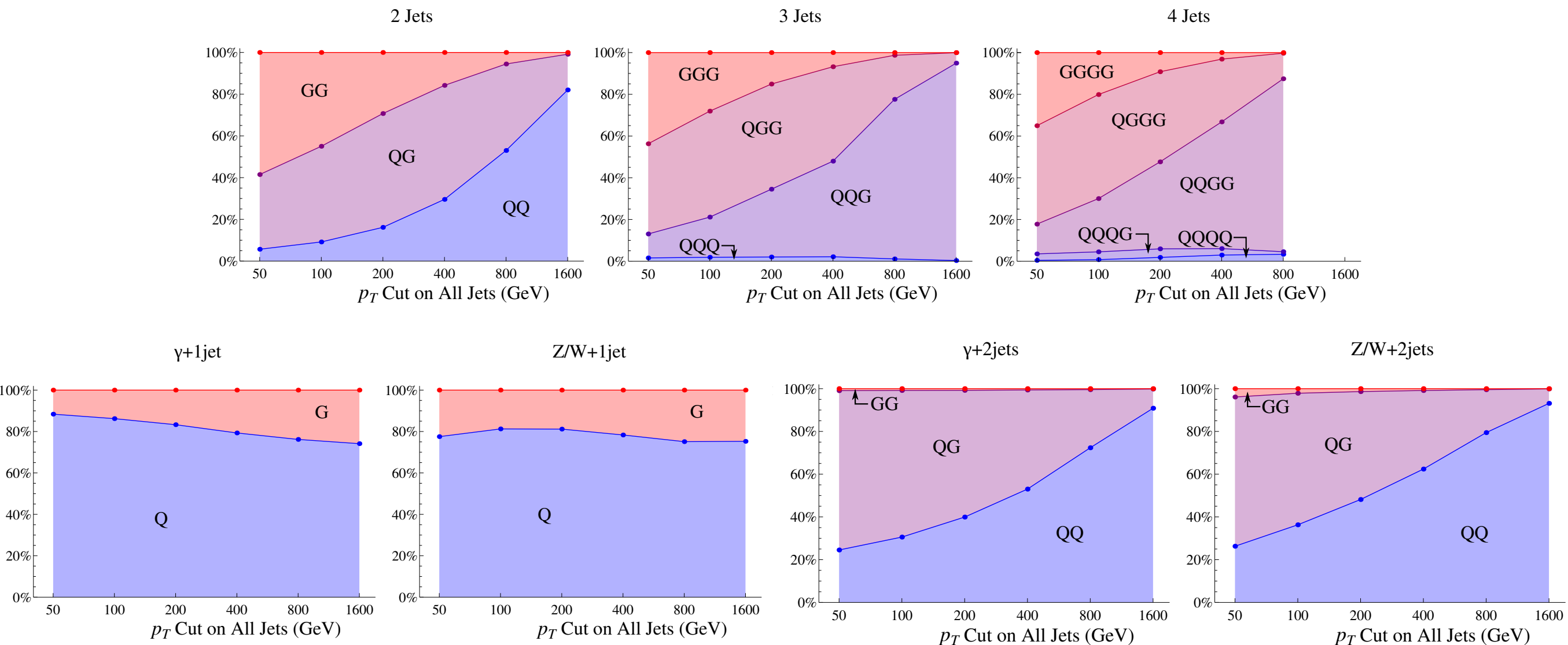Follow up work will use multiple templates

Apply to 3rd and 4th Jets

# Higher Jets Saw Larger Deviations

Transition from Quark Dominated Jets to Gluon Dominated Jets

Could hope to regress out the different compositions

Look at samples with different compositions

# Outlook

High Multiplicity Signals are Challenging
But Powerful Signal

$M_J$ & $N_J$ are powerful new tools to separate
new physics from QCD

Novel approaches to backgrounds exist
using Jet Factorization approximation

Learning how to have low background
searches without MET

# Thank You!

Boosted Community has been great to me

Grown from the small group in 2009

to this 115 person conference in its 6th iteration