



Measurement of quark and gluon jet fractions at the CMS: methods, results and outlook for Run-3

S. Shulha

JINR (RU)

F. Skorina GSU (BY)



XIXth Workshop on High Energy Spin Physics
dedicated to 90th anniversary of A.V. Efremov birth

4-8 September, 2023

JINR

- This work is part of the CMS analyses, which deals with **recognition** and **tagging** of q- and g-jets
- Recognition of q/g-jets is based on the **discriminator** - each jet is assigned a discriminator value V
- Examples of V are simple Macro Parameters (MP's): particle multiplicity inside jet (gluon jets have 1.5 times greater multiplicity), jet radius in (η, φ) -space (gluon jet is wider) or combination of simple MP's (**QGL** – “quark-gluon likelihood”,...)
- Discriminator is “trained” on MC jets: “**training**” means obtaining a MC normalized distributions over V for q/g-jets $\rightarrow H_{MC}^g(V)$ and $H_{MC}^q(V)$ –
 $H_{MC}^g(V)$ and $H_{MC}^q(V)$ are also called “**q/g-templates**”
- “q/g-templates” are **key objects** in q/g-tagging: “q/g-templates” allow one to say whether a given jet is a q- or g-jet with a given probability
- True “q/g-templates” $H_{DAT}^f(V)$ in data differ from model ones: $H_{DAT}^f(V) \neq H_{MC}^f(V)$
- Calculation of $H_{DATA}^f(V)$ using data is referred to as obtaining “**data-driven Scale Factor**” (**SF**) for q/g-templates: $S^f(V) \equiv H_{DAT}^f/H_{MC}^f$. SF is a **key issue** in q/g-tagging task

- To obtain SF (or q/g-templates) we need two jet samples with known g-fractions
- To date (Sept 2023), the official CMS recommendation for RUN-1 and RUN-2 is to use MC fractions for two channels (dijets and Z+jets) - $\alpha_{1,MC}^g$ and $\alpha_{2,MC}^g$:

$$H_{1,DAT} = \alpha_{1,MC}^g \cdot H_{DAT}^g + (1 - \alpha_{1,MC}^g) \cdot H_{DAT}^q \quad (1)$$

$$H_{2,DAT} = \alpha_{2,MC}^g \cdot H_{DAT}^g + (1 - \alpha_{2,MC}^g) \cdot H_{DAT}^q$$

- Solution of this system of Eqs. gives us data-driven corrected **q/g-templates**:

$$H_{DAT}^q = \frac{\alpha_{2,MC}^g H_{1,DAT} - \alpha_{1,MC}^g H_{2,DAT}}{\alpha_{2,MC}^g - \alpha_{1,MC}^g} \quad (2)$$

$$H_{DAT}^g = (g \rightarrow q, 1 \leftrightarrow 2)$$

- 1st **recommendation for us** was to apply SF in measurement of g-fraction
- But, in current official form, Eqs.(2) were written w/o normalization and with **hidden** MC g-fractions. It is not difficult to guess from Eqs.(1) and (2) that measured g-fraction with corrected q/g-templates in the data will give **exactly** the MC g-fractions!

Tip for the careful listener: measured $\alpha_{1,DAT}^g$ is a solution of Eq. like 1st Eqs (1):

$$H_{1,DAT} = \alpha_{1,DAT}^g \cdot H_{DAT}^g + (1 - \alpha_{1,DAT}^g) \cdot H_{DAT}^q \quad (1')$$

$$\alpha_{1,DAT}^g = \alpha_{1,MC}^g$$

- **We proposed (2020)** to use in CMS the modified SF for q/g-templates:

$$H_{\text{DAT}}^q = \frac{\alpha_{2,\text{DAT}}^g H_{1,\text{DAT}} - \alpha_{1,\text{DAT}}^g H_{2,\text{DAT}}}{\alpha_{2,\text{DAT}}^g - \alpha_{1,\text{DAT}}^g} \quad (3)$$

$$H_{\text{DAT}}^g = (q \leftrightarrow g, 1 \leftrightarrow 2)$$

- Before obtaining SF and $H_{\text{DAT}}^{q/g}(V)$ we **need to measure g-jet fractions**. So, measurement of g-jet fraction becomes a **key task** for q/g-tagging!

We have found another important correction to SF (3):

- Eqs.(3) give universal q/g-templates for any channel and any jet kinematics. But, MC q/g-templates depend on kinematics! **We proposed** method to introduce in Eqs.(3) **kinematical non-universal terms** (SS, D.Budkouski, PEPAN Lett 2021-2022)

Very important remark:

- g-fraction measurement with corrected q/g-templates Eqs.(3) gives the same $\alpha_{1,\text{DAT}}^g$ and $\alpha_{2,\text{DAT}}^g$. So, 1st measurement of g-fractions with MC q/g-templates **cannot be improved by SF** – iteration process is impossible!

Proposition: $\alpha_{1,\text{DAT}}^{g'} \equiv \alpha_{1,\text{DAT}}^g$

Tip for the careful listener: to prove this, we need to write two equations

1st iteration $\alpha_{1,\text{DAT}}^g$ is a solution of Eq.: $H_{1,\text{DAT}} = \alpha_{1,\text{DAT}}^g \cdot H_{\text{MC}}^g + (1 - \alpha_{1,\text{DAT}}^g) \cdot H_{\text{MC}}^q$

2nd iteration $\alpha_{1,\text{DAT}}^{g'}$ is a solution of Eq.: $H_{1,\text{DAT}} = \alpha_{1,\text{DAT}}^{g'} \cdot H_{\text{DAT}}^g + (1 - \alpha_{1,\text{DAT}}^{g'}) \cdot H_{\text{DAT}}^q$

Tip for the careful listener (cont.):

$$H_{1,\text{DAT}} = \alpha_{1,\text{DAT}}^{g'} \cdot H_{\text{DAT}}^g + (1 - \alpha_{1,\text{DAT}}^{g'}) \cdot H_{\text{DAT}}^q$$

$$H_{\text{DAT}}^q = \frac{\alpha_{2,\text{DAT}}^g H_{1,\text{DAT}} - \alpha_{1,\text{DAT}}^g H_{2,\text{DAT}}}{\alpha_{2,\text{DAT}}^g - \alpha_{1,\text{DAT}}^g}$$

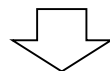
$$H_{\text{DAT}}^g = \frac{(1 - \alpha_{1,\text{DAT}}^g) H_{2,\text{DAT}} - (1 - \alpha_{2,\text{DAT}}^g) H_{1,\text{DAT}}}{\alpha_{2,\text{DAT}}^g - \alpha_{1,\text{DAT}}^g}$$

$$\alpha_{1,\text{DAT}}^{g'} = \frac{H_{1,\text{DAT}} - H_{\text{DAT}}^q}{H_{\text{DAT}}^g - H_{\text{DAT}}^q}$$

$$H_{1,\text{DAT}} - H_{\text{DAT}}^q = \frac{\alpha_{1,\text{DAT}}^g (H_{2,\text{DAT}} - H_{1,\text{DAT}})}{\alpha_{2,\text{DAT}}^g - \alpha_{1,\text{DAT}}^g}$$

$$H_{\text{DAT}}^g - H_{\text{DAT}}^q = \frac{H_{2,\text{DAT}} - H_{1,\text{DAT}}}{\alpha_{2,\text{DAT}}^g - \alpha_{1,\text{DAT}}^g}$$

$$\alpha_{1,\text{DAT}}^g = \frac{H_{1,\text{DAT}} - H_{\text{MC}}^q}{H_{\text{MC}}^g - H_{\text{MC}}^q} \quad (4)$$



Proposition : $\alpha_{1,\text{DAT}}^{g'} \equiv \alpha_{1,\text{DAT}}^g \quad \otimes$

- **2nd iteration** for g-fraction measurement is **impossible!**
- The model determines g-fraction in the data unambiguously and does not allow it to be corrected if the original model is not changed
- However, there is a way to define **quantitatively** discrepancy between model q/g-templates and data ones in measured g-fractions – it is **Model Uncertainty** (M.U.)
- To find M.U., we need to use **several** independent jet **macro parameters**...

- If $\alpha_{\text{DAT}}^g \approx \alpha_{\text{MC}}^g$ then official SF \approx new SF
- Spoiler: we found **strong g-jet suppression** in region $P_T^{\text{jet}} < 200$ GeV:

$$\alpha_{\text{DAT}}^g \approx (0.5 \div 0.7) \cdot \alpha_{\text{MC}}^g \quad \Rightarrow \quad \text{official SF} \gg \text{new SF}$$

- Thus, official CMS SF's developed for Run-1 and Run-2 are wrong: they correct the g-factions $\alpha_{\text{DAT}}^g \rightarrow \alpha_{\text{MC}}^g$

It is our negative contribution to CMS "q/g-tagging".
It should be taken into account in CMS Run-3 analyses

Now we are moving to g-fraction measurements...

- Careful listener may suggest a **method** for measuring – the main formula has already been written on page 5:

$$\alpha_{\text{DAT}}^g = \frac{H_{\text{DAT}} - H_{\text{MC}}^q}{H_{\text{MC}}^g - H_{\text{MC}}^q} \quad (4)$$

where $H_{\text{DAT}}(V)$ – measured distribution, $H_{\text{MC}}^f(V)$ - MC q/g-templates

- But right part depends on V -bin?
- Well! **Each V -bin** can be considered as **independent experiment** and we define measured α_{DAT}^g as averaged value...

Method of "bin averaging"

S.S. PEPAN Lett. **2023/2024** (in preparation)

- For any MP (jet macro parameter) $V \equiv V_{1,2,3,4,\dots}$:

$$H^{\text{MC,DAT}}(V) = \alpha^g H^g(V) + (1 - \alpha^g) H^q(V) \quad (5)$$

In case of MC, Eq.(5) has the same solution α^g for all V -bins:

$$\alpha^g = \frac{H^{\text{MC}}(V) - H^q(V)}{H^g(V) - H^q(V)} = \text{const}(V)$$

- In case of DATA, solution of Eq. (5) is not a V -constant:

$$\alpha_V^g = \frac{H^{\text{DATA}}(V) - H^q(V)}{H^g(V) - H^q(V)} \quad (6)$$

Each bin is a separate independent experiment to measure α^g

- Definition:** measured g-fraction is averaged value:

$$\alpha^g \equiv \langle \alpha_V^g \rangle = \frac{\sum_{V=1}^{N_V} \alpha_V^g}{N_V} \quad (7)$$

N_V - number of V -bins

$$\text{with uncertainty } \Delta\alpha^g \equiv \frac{\sqrt{\langle \alpha_V^{g2} \rangle - \langle \alpha_V^g \rangle^2}}{\sqrt{N_V}}$$

- In **June 2023** we implemented this method and showed results in CMS
- Deprecated method:** So far, we have used a more complex method with QGL and with fit:

$$\text{WLS or LS methods by ROOT/MINUIT: } H_{\text{DAT}} \sim \alpha_{\text{DAT}}^g \cdot H_{\text{MC}}^g + (1 - \alpha_{1,\text{DAT}}^g) \cdot H_{\text{MC}}^q$$

- Method of "bin averaging" allows to find stable result for **any/all jet MP's** $V_{1,2,3,4,\dots}$ with small statistics

- In case of MC, calculation with any jet MP $V_{1,2,3,\dots}$ gives the same $\alpha_1^g = \alpha_2^g = \alpha_3^g = \dots = \alpha^g$ because **q/g-templates are true for MC**

$$\alpha^g = \frac{H(V_k) - H^q(V_k)}{H^g(V_k) - H^q(V_k)} = \mathbf{const(k)}$$

- In case of DATA, calculation with any MP $V_{1,2,\dots}$ gives different $\alpha_1^g \neq \alpha_2^g \neq \alpha_3^g \neq \dots$ because **MC q/g-templates are not true for DATA**

- Maximum of differences $|\alpha_1^g - \alpha_2^g|, |\alpha_1^g - \alpha_3^g|, |\alpha_2^g - \alpha_3^g|, \dots$ describes the deviation of MC q/g-templates from true ones = **Model Uncertainty (M.U.)**

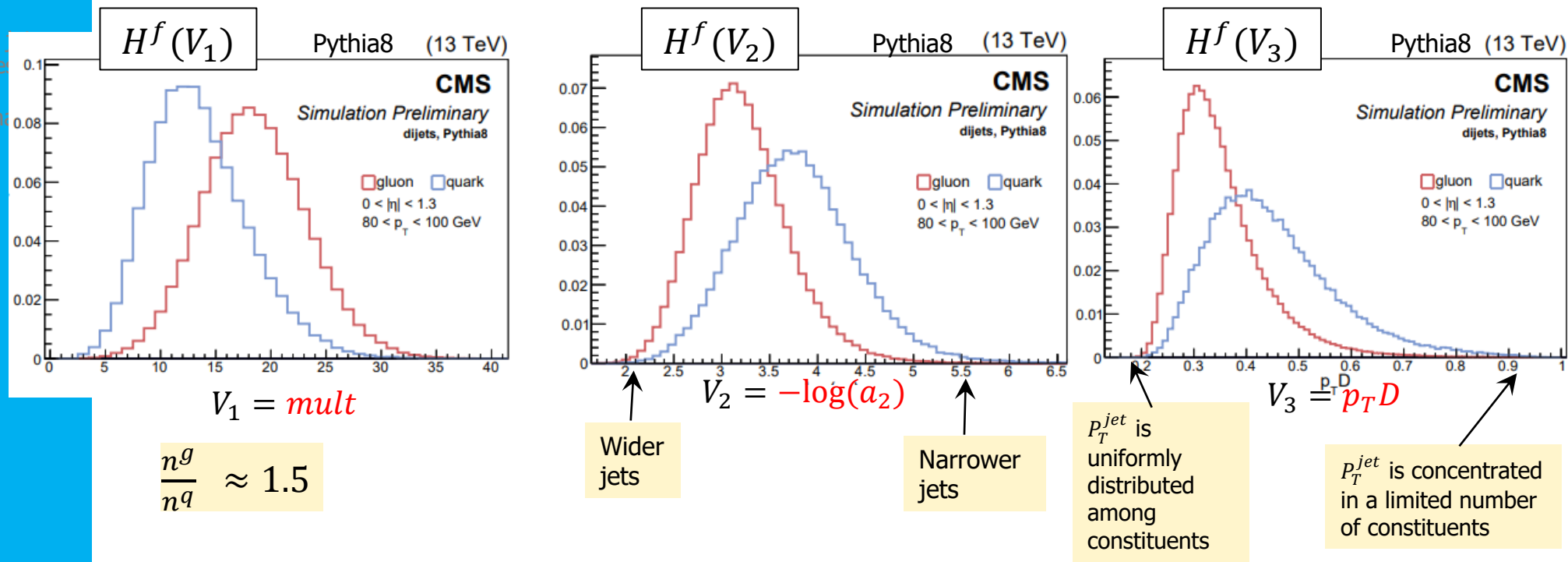
$$\mathbf{M.U.} = \frac{1}{2} \cdot \max\{|\alpha_1^g - \alpha_2^g|, |\alpha_1^g - \alpha_3^g|, |\alpha_2^g - \alpha_3^g|, \dots\}$$

Choose MP's which are the most sensitive to Jet Flavour¹

- Total multiplicity inside jet (*mult*)
- Minor axis of jet ellipse in (η, ϕ) -space a_2
- "Fragmentation function" $p_T D = \frac{\sqrt{\sum_i p_{T i}^2}}{\sum_i p_{T i}} \in [0, 1]$

$$V_{1,2,3} = (\text{mult}, a_2, p_T D) \equiv \vec{V}$$

Fig. 1: q/g-templates $H^f(V_1), H^f(V_2), H^f(V_3)$



These three q/g-templates are used to measure g-fractions

- QGL is a jet MP that is a combination of simple MP's :

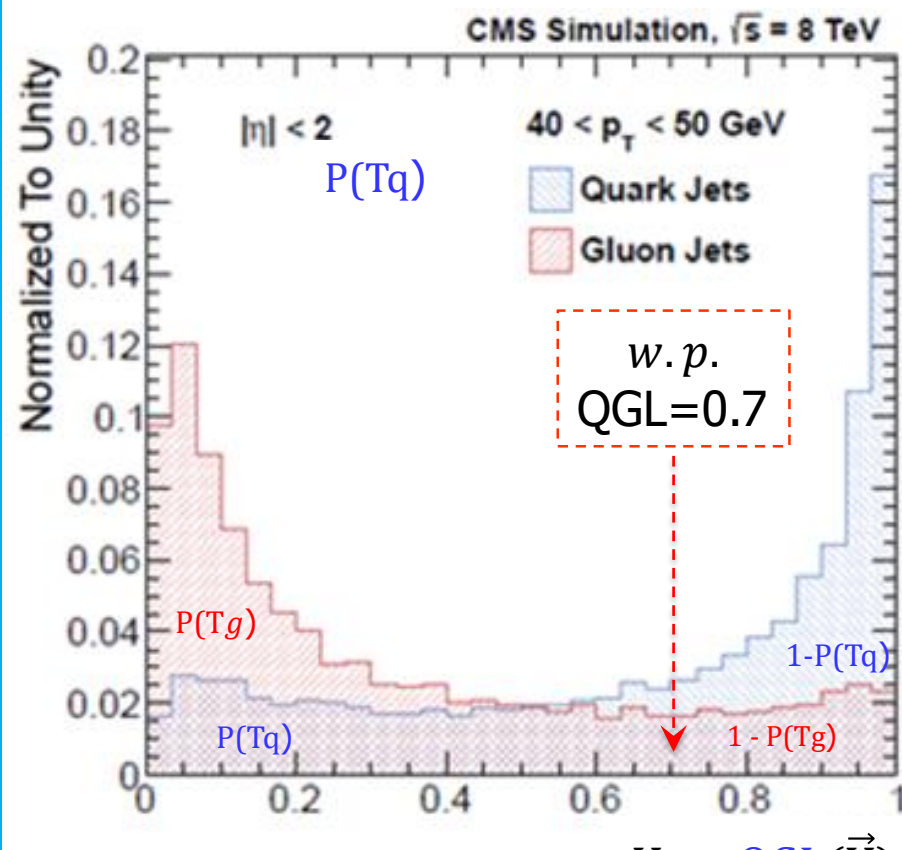
$$V_4 \equiv QGL = \frac{Q(\vec{V})}{Q(\vec{V}) + G(\vec{V})}$$

QGL – discriminator
 "Quark-Gluon Likelihood"

$$Q(\vec{V}) = \prod_{i=1}^3 H^q(V_i), \quad G(\vec{V}) = \prod_{i=1}^3 H^g(V_i)$$

$$V_{1,2,3} = (mult, a_2, p_T D) \equiv \vec{V}$$

- Sensitivity of QGL to jet flavour is much stronger than that of original $mult, a_2, p_T D$.



With respect to w.p.:

- Left/right red area= g-jet efficiency/rejection
- Left/right blue area= q-jet rejection/efficiency

- QGL-templates are used to tag q/g-jets. It is very important tool to select channels
- We measured g-fractions with QGL-templates to **check QGL** written in datasets

- We show (**July 2023**) that QGL written in **all CMS Run-2 datasets** are **wrong**
- We prepared **new QGL** for CMS Run-2 and test them using g-fraction measurements

It is our contribution to CMS "q/g-tagging"
 It should be implemented in Run-2(2016-2018) analyses

Fig.2: QGL-templates

$$V_4 \equiv QGL(\vec{V})$$

- α^g was found by $V = mult, a_2, p_T D$ and "new QGL"

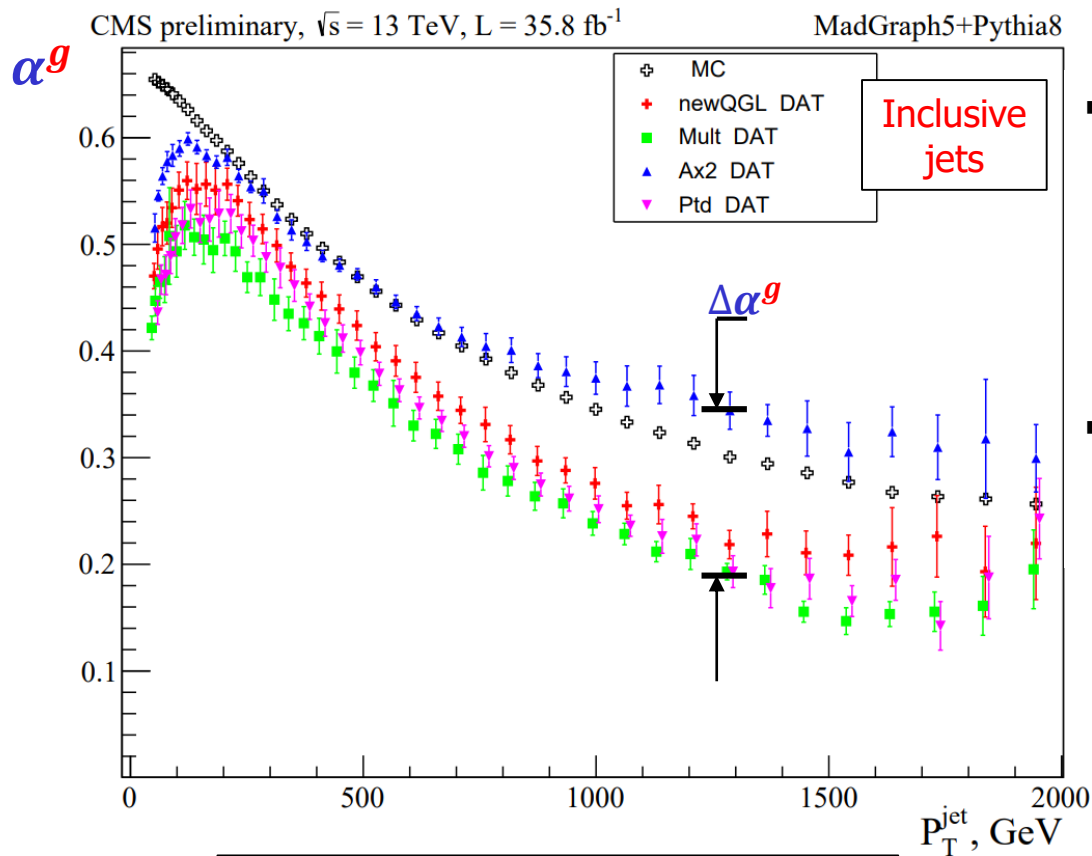


Fig. 3: Demonstration of M.U.

- Measurement of g-fraction demonstrates indirectly large deviation of true unknown DATA q/g-templates from Pythia8 ones

- This preliminary results were obtained in CMS group "Gluon-jet/Quark-jet analyses" ¹:

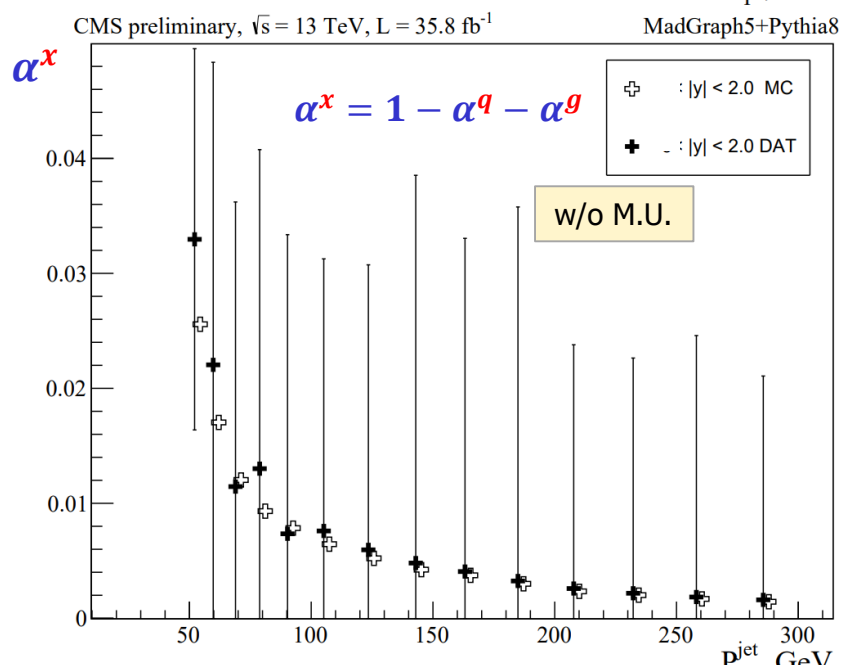
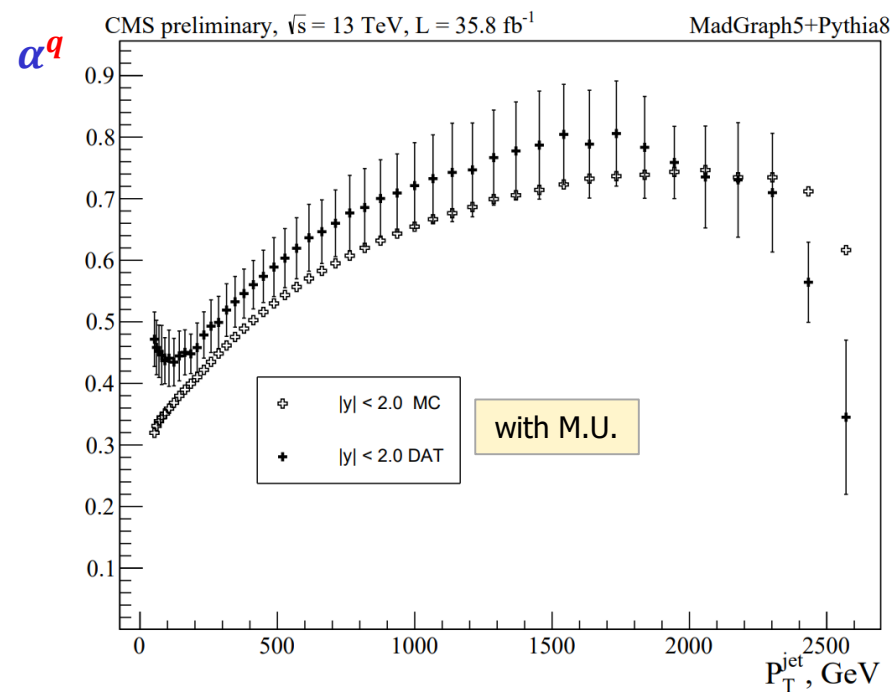
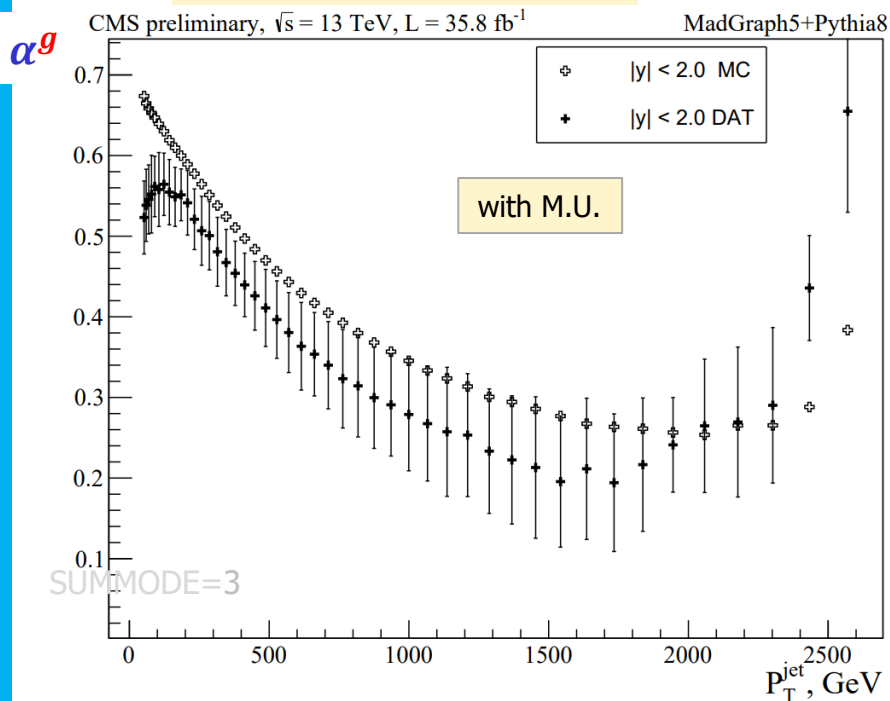
S.S., D.Budkouski(JINR), J.Strologas (GR), O.Atakisi(TR)

- This group was created in April 2021 purposefully to measure g-fractions in inclusive jet channel with Run-II data

¹<https://indico.cern.ch/category/12755/>

q/g/x-jet fractions

Inclusive jets



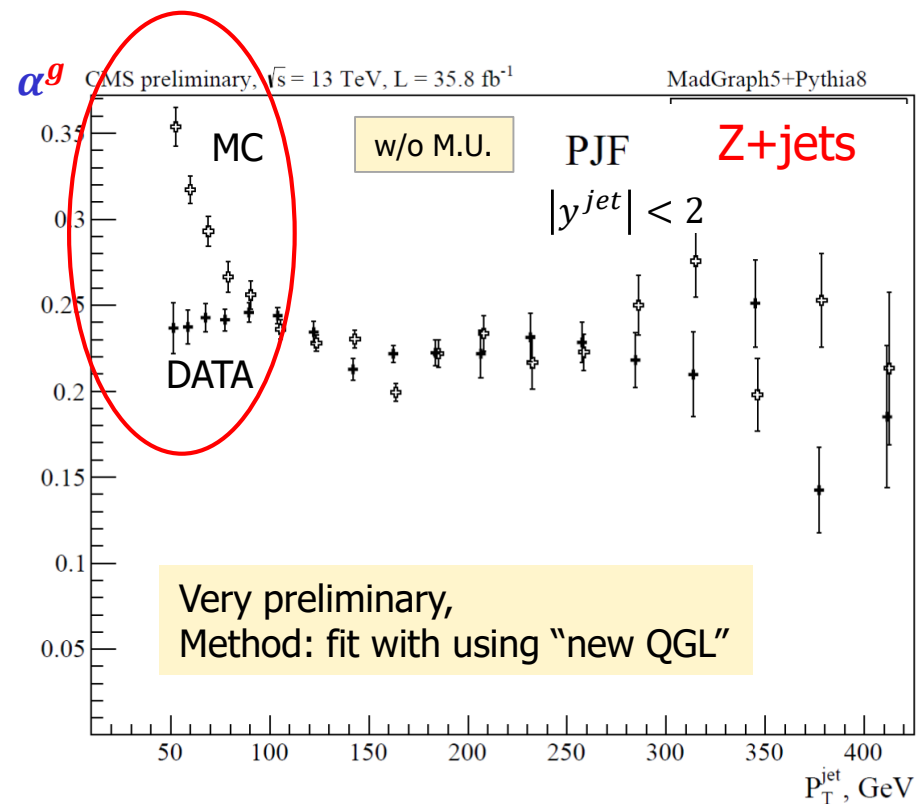
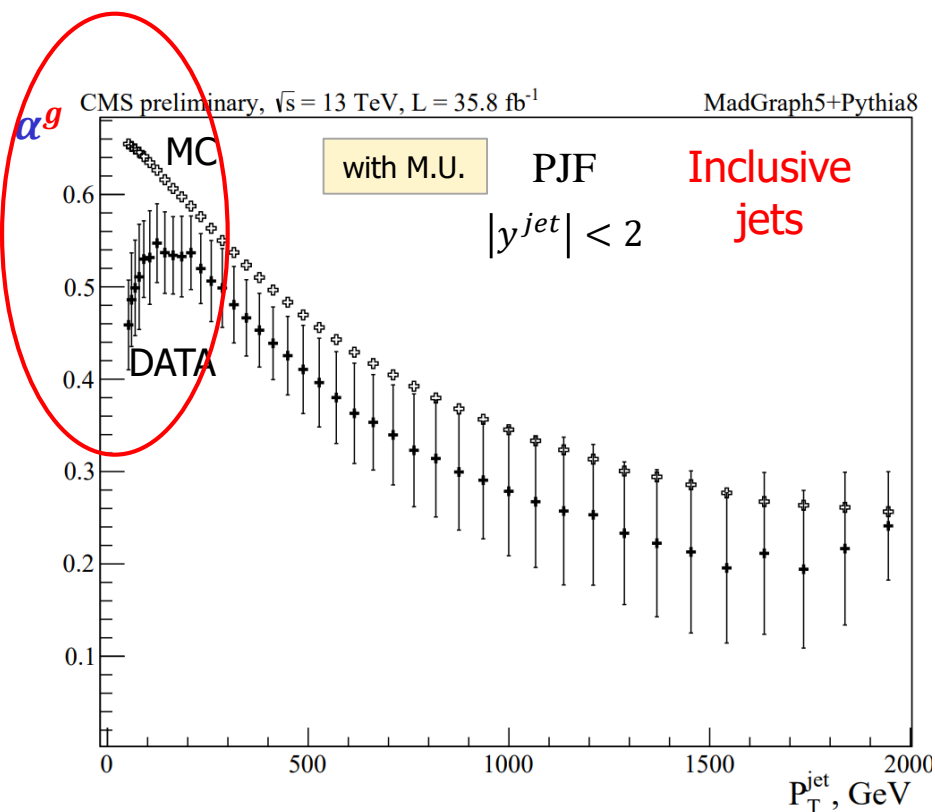
- This measurement results were obtained in CMS group "Gluon-jet/Quark-jet analyses" ¹:

S.S., D.Budkouski(JINR), J.Strologas (GR), O.Atakisi(TR)

¹<https://indico.cern.ch/category/12755/>

Run-II(2016)

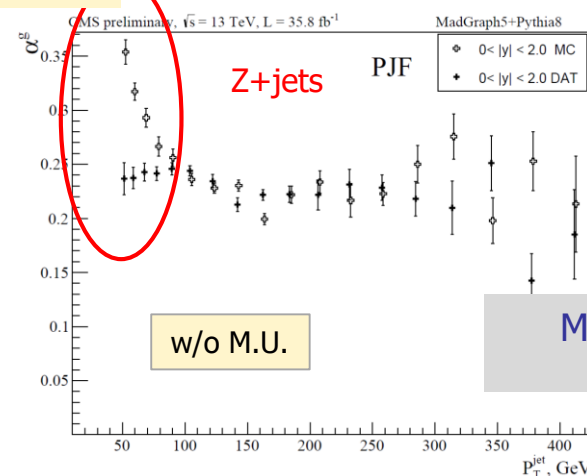
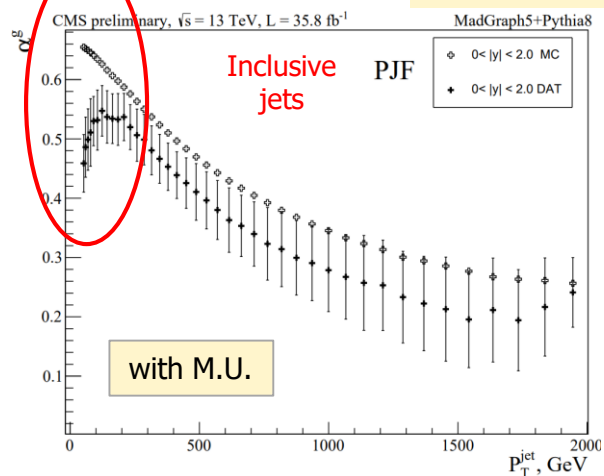
- α^g -jet suppression is visible at low P_T^{jet} in "Inclusive jets" and in "Z+jets"



MadGraph5+Pythia8

ak4-jets: $R = 0.4$

Run-II(2016)



- Similar results we obtained earlier for **Run-I (2012)**
- **Run-I** results are documented:

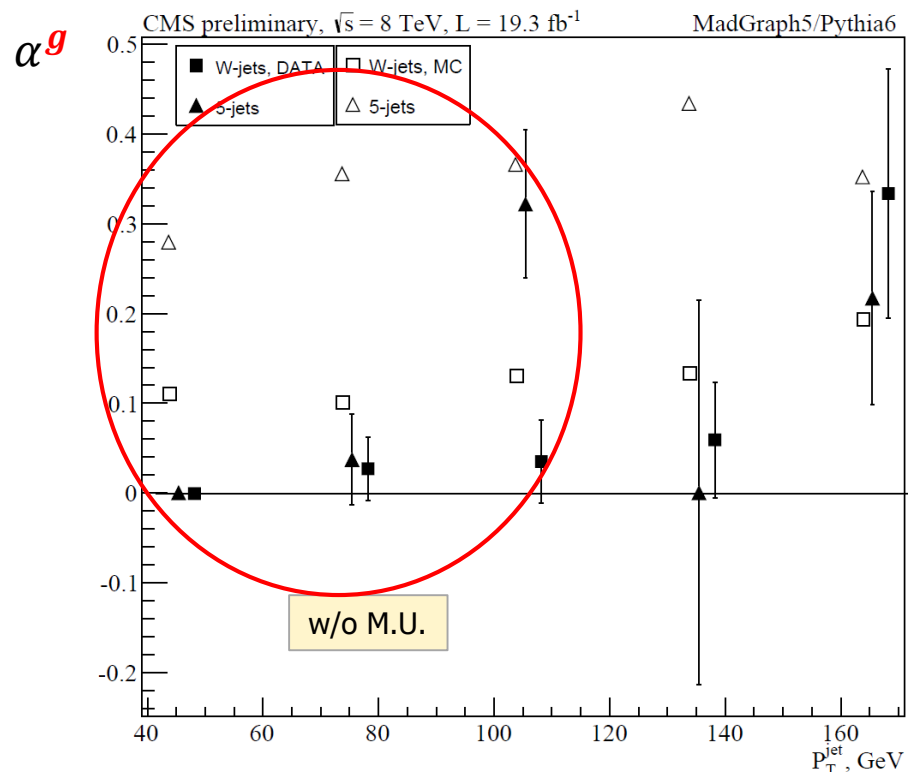
S.S., S.Shmatov, A.Zarubin: CMS AN-2018-131, **2018**

S.S. D.Budkouski, CMS AN-2020-143, **2020**

S.S. D.Budkouski, CMS AN-2021-024, **2021**

S.S. SMP-HAD Workshop, 11 Feb **2020**, <https://indico.cern.ch/event/861896/>

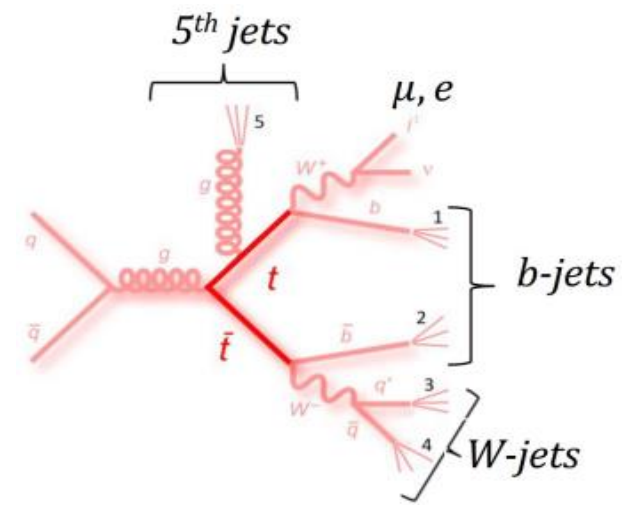
S.S. SMP-HAD Meeting, 1 June **2018**, <https://indico.cern.ch/event/732652/>



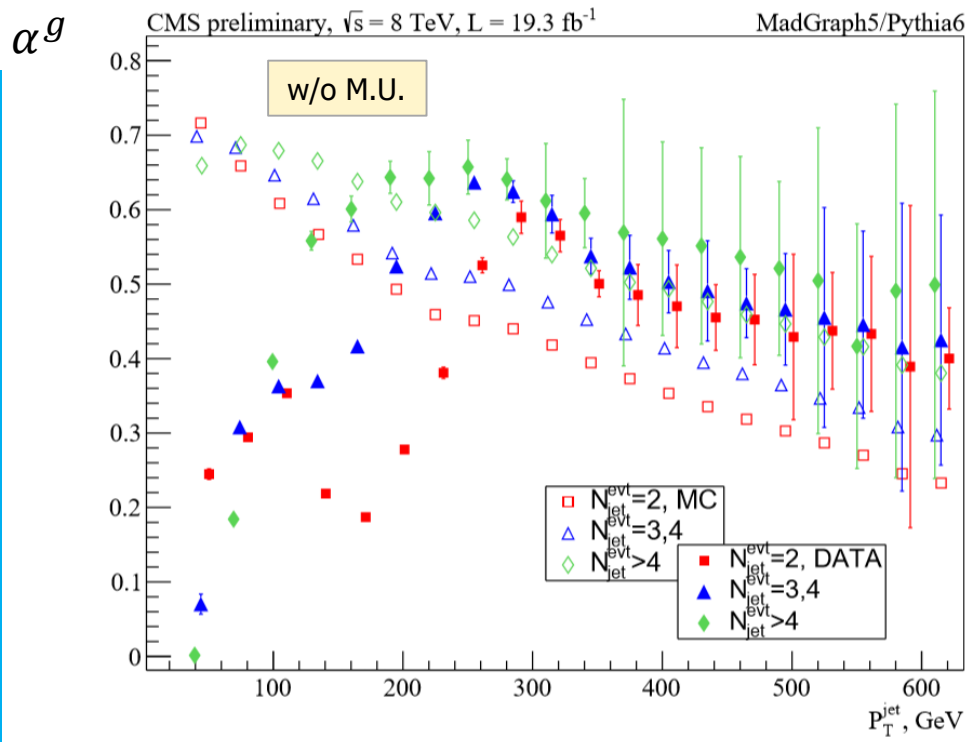
MadGraph5+Pythia6

ak5-jets: **R = 0.5**

- Semileptonic $t\bar{t}$ channel
- **M.U.** is not shown



| N_{jets} | Jet name | $P_T^{jet}, \text{ GeV}$ | $\alpha_k^{g,DATA}, \%$ | $\alpha_k^{g,MC}, \%$ |
|------------|-----------------------|--------------------------|-------------------------|-----------------------|
| 4 | W-jets | 30 ÷ 150 | 0 ÷ 5 (± 5) | 10 ÷ 11 |
| ≥ 5 | 5 th -jets | 30 ÷ 90 | 0 ÷ 3 (± 5) | 28 ÷ 34 |



MadGraph5+**Pythia6**
 ak5-jets: **R = 0.5**

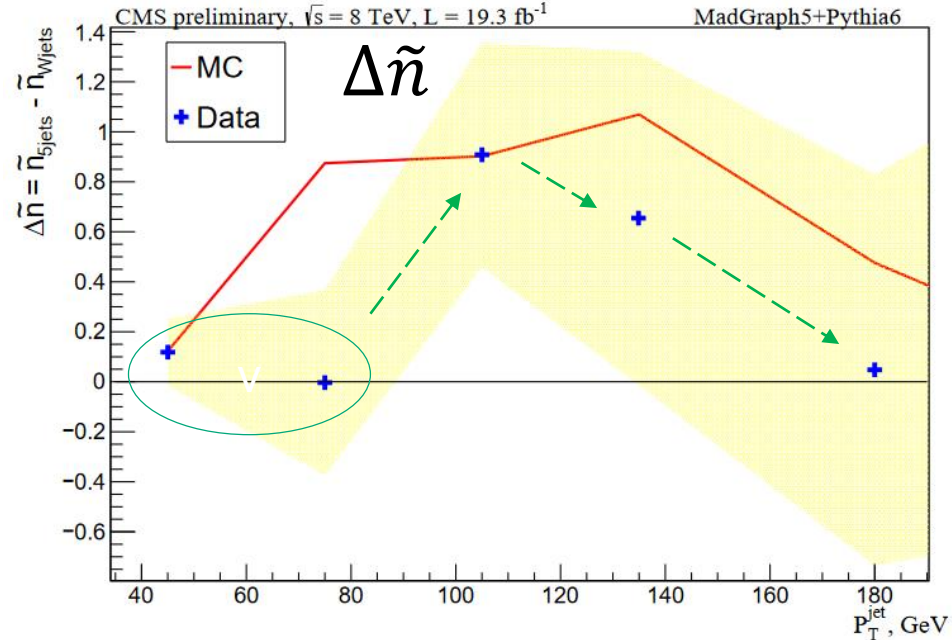
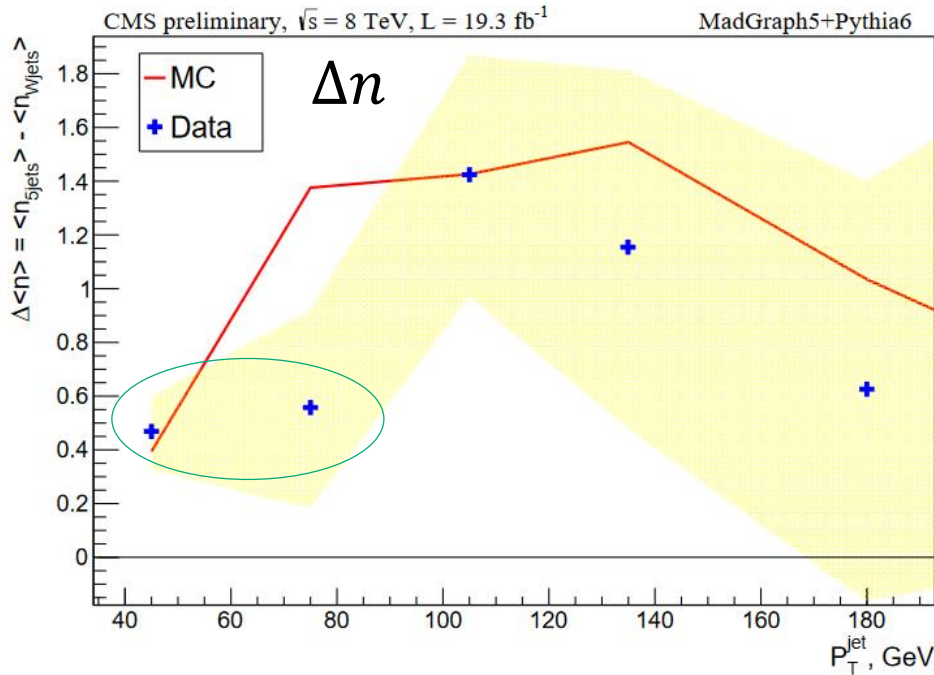
- Dijet, Run-I(2012)
- **M.U.** is not taken into account
- HLT prescaling is not taken into account

Name

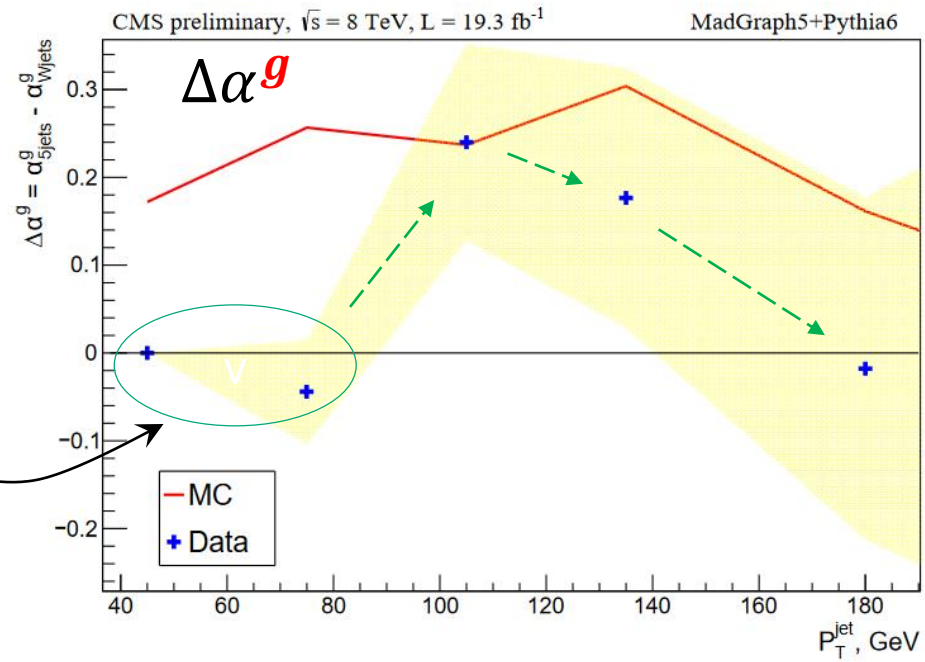
| N_{jets}^{evt} | P_T^{jet} , GeV | $\alpha_k^{g,DATA}$, % | $\alpha_k^{g,MC}$, % | Name | |
|------------------|-------------------|-------------------------|-----------------------|-------------------|--------------------------|
| 2 | 30 ÷ 210 | 16 ÷ 35 | 72 ÷ 50 | “dijet-1” (red) | |
| 3,4 | 30 ÷ 180 | 6 ÷ 40 | 70 ÷ 60 | “dijet-2” (blue) | |
| ≥ 5 | 30 ÷ 120 | 0 ÷ 40 | 65 ÷ 69 | “dijet-3” (green) | |
| 4 | 30 ÷ 150 | 0 ÷ 5 (± 5) | 10 ÷ 11 | W-jets | Semi-leptonic $t\bar{t}$ |
| ≥ 5 | 30 ÷ 90 | 0 ÷ 3 (± 5) | 28 ÷ 34 | 5th-jets | |

Run-I(2012) semileptonic $t\bar{t}$

$$A \cdot \Delta\tilde{n} = \Delta\alpha^g$$



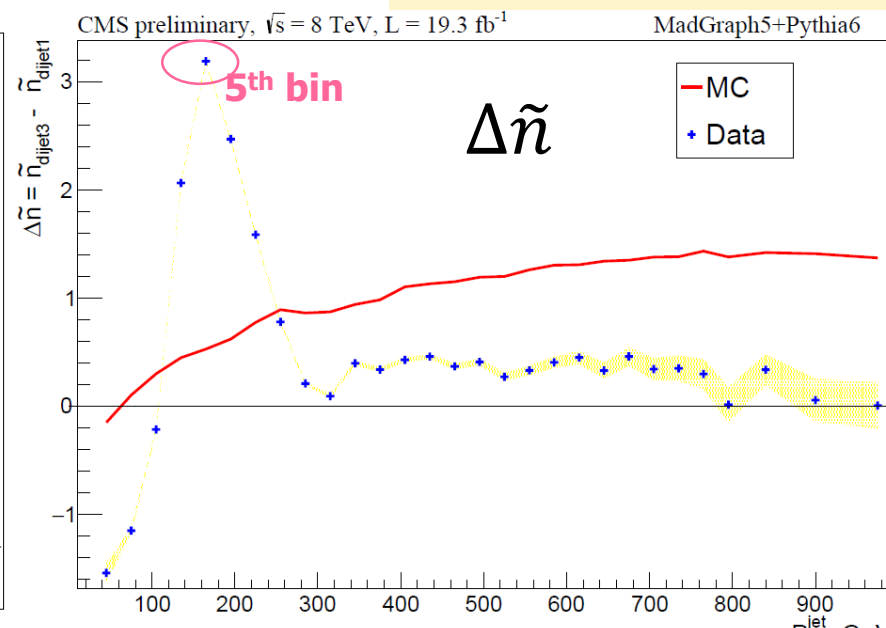
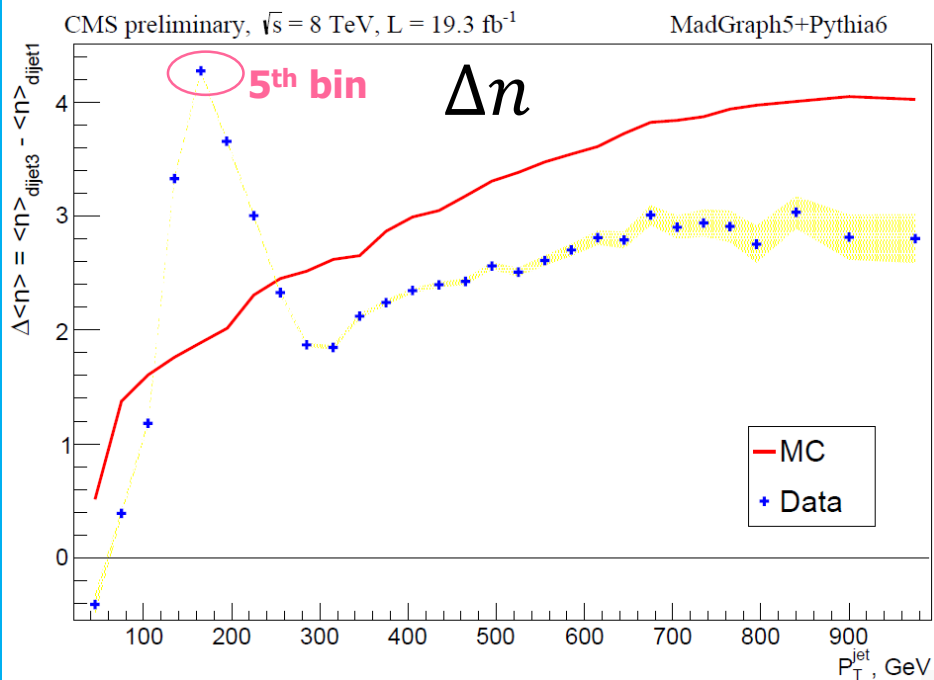
- $\Delta\tilde{n}$ and $\Delta\alpha^g$ are similar:
 $\Delta\tilde{n} = A \Delta\alpha^g \approx 0$ in 1st and 2nd bins !
- Measurement of mean jet **C.P.M**'s indirectly confirms g -jet suppression



"Test"

Run-I(2012) Dijet

$$A \cdot \Delta\tilde{n} = \Delta\alpha^g$$

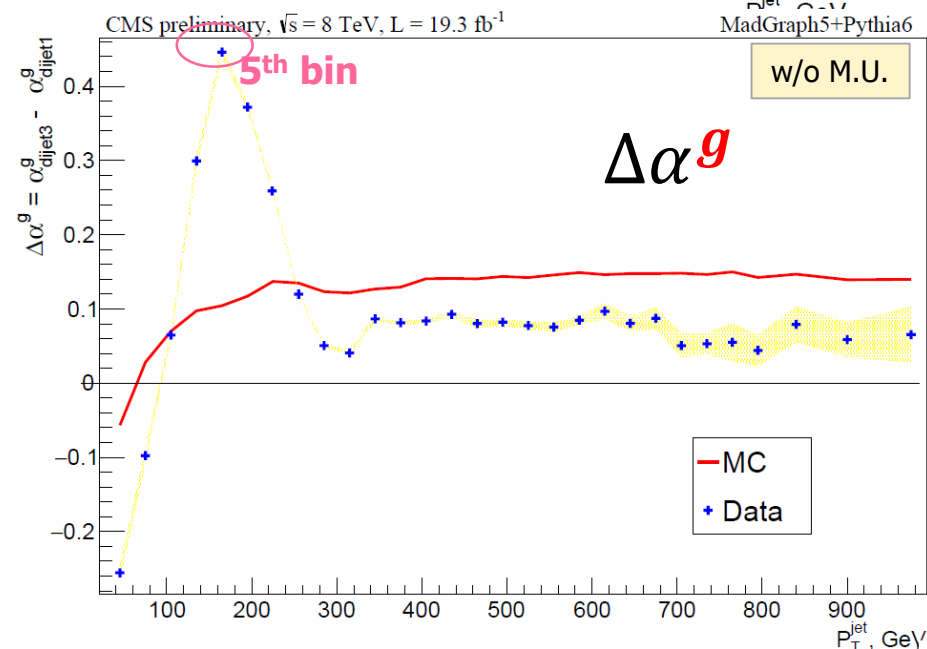


- $\Delta\tilde{n}$ and $\Delta\alpha^g$ are **similar** in all bins:

$$A \cdot \Delta\tilde{n} = \Delta\alpha^g \quad !$$

- Measurement of mean jet **C.P.M's** indirectly **confirms g -jet suppression at low P_T^{jet}**

"Test"

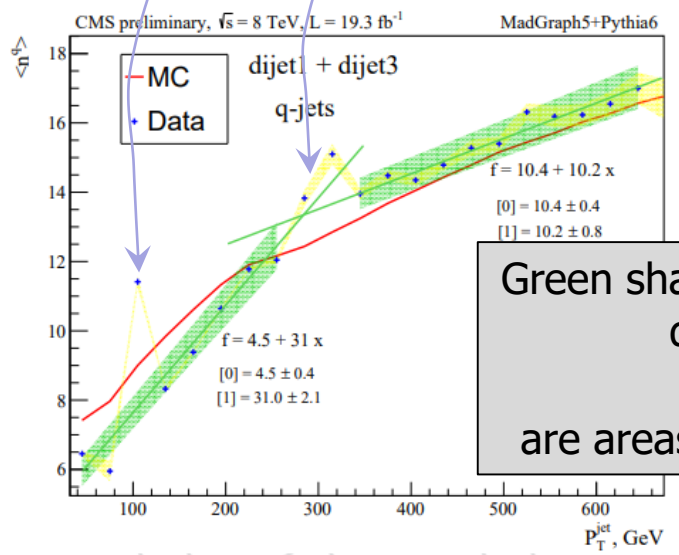
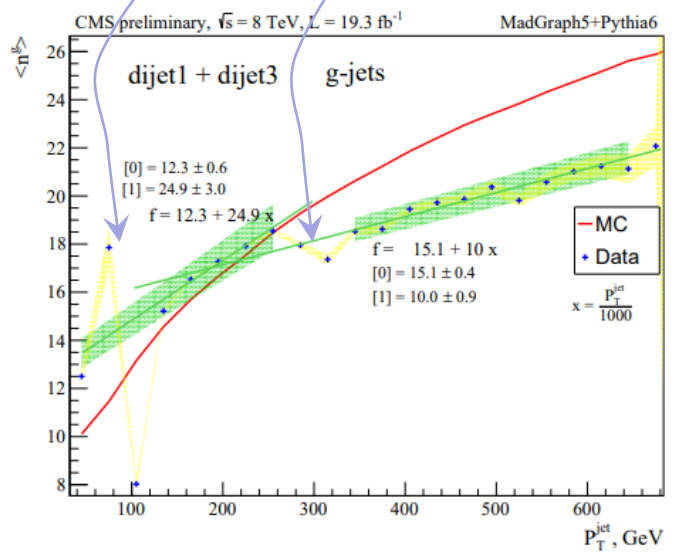


- Measurement of g-fractions in many channels was proposed, developed and implemented in CMS (Run-1 and Run-2)
- It was shown that g-fraction measurement is a 1st stage in preparation of q/g-templates used in q/g-tagging
- Possible phenomenon of g-jet suppression in low P_T^{jet} region is observed by all studied channels, for CMS Run-1 and Run-2

Run-I(2012) Dijet

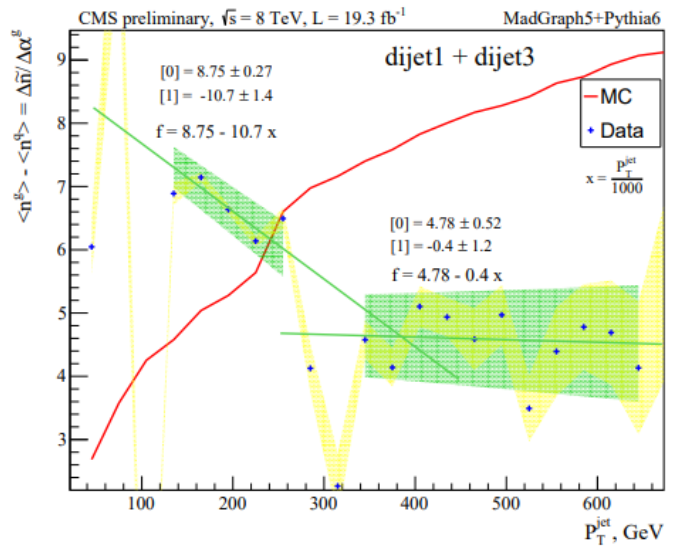
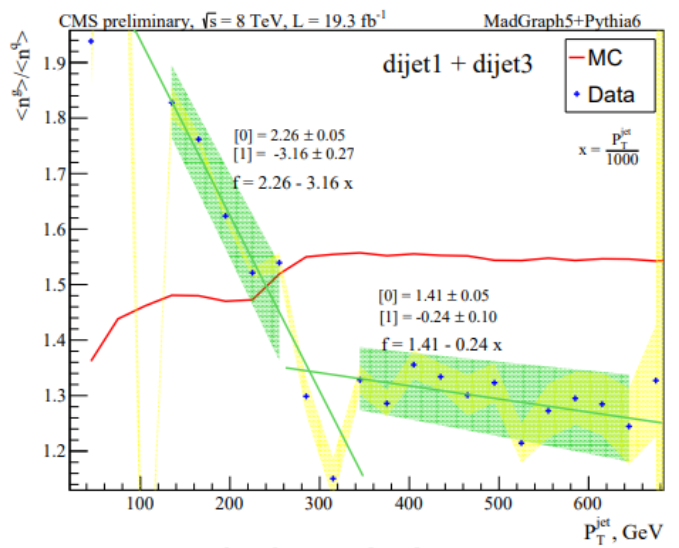
w/o M.U.

1st measurement of q/g-jet mean C.P.M's using measured g-fractions in two channels



Small denominator $\alpha_1^g - \alpha_2^g \sim 0$ in formulas makes solutions unstable

Green shaded areas with large denominator $\alpha_1^g - \alpha_2^g$ are areas of robust solutions



Two dijet samples used:
 dijet1: $N_{jets}^{evt} = 2$
 dijet3: $N_{jets}^{evt} \geq 5$