



# CEPC Workshop Status and Outlook of H->bb/cc/gg Analysis

Yu Bai (from Southeast University, Nanjing) On Behalf of CEPC Physics-Software Study Group November 19, 2019



## Why H->bb/cc/gg is important?



A measurement on the couplings between Higgs and quarks will decode the origin of mass in great confidence

### **Current Status**

H->bb/cc/gg is expected to be 57%/9%/3%



#### NOT bad with a 125 GeV Higgs mass

But that's not an easy task for our LHC colleges!

#### **Direct H->bb measurement**

- Observation of H->bb and VH production by ATLAS Collaboration, <u>Phys.Lett.B 786 (2018) 59</u>
- Observation of Higgs boson decay to bottom quarks by CMS Collaboration, <u>Phys.Rev.Lett (2018) 121801</u>

The H->bb signal strength was measured with precision around 20%, consisted with SM prediction

**Direct H->cc measurement:** 

 Search for the decay of the Higgs boson to Charm Quarks with the ATLAS Experiment Phys. Rev. Lett (2018), 211802

A 95% CL upper limit set at about signal strength = 100

#### **Gluon-gluon fusion Analysis**

- <u>H->ττ in ATLAS</u>
- <u>H->WW in ATLAS</u>
- <u>H-> *ττ* in CMS</u>
- <u>H->γγ in CMS</u>

Uncertainty of O(10%)

# Review of Previous Study

Study based on pre-CDR set up, using full simulated sample with sqrt(s)=250 GeV and old geometry (some extrapolate to sqrt(s) = 240 GeV)

## H->bb/cc/gg at CEPC



# Analysis Strategy



#### **Dominant Backgrounds include the ZZ/WW events with same final states**

### **Event Selection**





#### recoil mass is the crucial variable to distinguish signal and background

	signal	Signal eff	hiqas bka	non-hiqqs bkq
ee channel	9.15k	52.6%	1.10k	6.15k
<i>uu</i> channel	12.8k	63.9%	1.48k	5.29k
$\nu\nu$ H				



Hpp



Mass  $\chi^2$  with different hypothesis are combined to reject ZZ/WW events

#### qqΗ

signal	Sign al eff	higgs bkg	4f-hadronic	qq	4f- semilept
211.2k	42.8%	32.6k	1.08M	405.6k	0.58k

#### $\nu\nu$ H

signal	Signal eff	higgs bkg	non-higgs bkg
85.8k	49.2%	1.96k	22.88k

7

# Flavor Tagging and Flavor template

#### **Flavor tagging performance**









- 4 categories of events according to the soft-leptons and vertices multiplicity
- GBDT with tagging-sensitive variables method applied to each category

#### The performance of FT ensures the goal of precision

Processes with different flavor components can be separated by template fit



## Template-Recoil Mass Combined Fit

#### Assuming lepton pair's recoil mass and jet flavor are independent in signal

 $\mathrm{PDF}^{3D}(X_B, X_C, M_{\mathrm{recoil}}) = \mathrm{PDF}^{flavor}(X_B, X_C) \times \mathrm{PDF}^{\mathrm{recoil\_mass}}(M_{\mathrm{recoil}})$ 



- The shape parameter of recoil mass in signal and dominate background are float in the fit
- Reduce the dependency to the MC prediction
- Effects of systematic uncertainty also considered

#### **Current Results**

**Combination of the 4 channels:** 

Statistic precision of  $\sigma$ (ZH)\*Br(H->bb/cc/gg) is 0.3% 3.3% and 1.3%

	Decay mode	$\sigma(ZH) \times BR$
Consistent with the goal expected	$H \rightarrow b\bar{b}$	0.28%
	$H \to c\bar{c}$	2.2%
in pre-CDR with full simulation samples	$H \to gg$	1.6%

IIH with 3D fit and systematic uncertainties considered:

				-		
	$\mu^+\mu^-H$			$e^+e^-H$		
	$H \rightarrow b\bar{b}$	$H \rightarrow c \bar{c}$	$H \rightarrow gg$	$H \rightarrow b \bar{b}$	$H \mathop{\rightarrow} c \bar{c}$	$H \mathop{\rightarrow} gg$
Statistic Uncertainty	1.1%	10.5%	5.4%	1.6%	14.7%	10.5%
Fixed Background	-0.2%	+4.1%	7.6%	-0.2%	+4.1%	7.6%
Tixed Dackground	+0.1%	-4.2%	1.070	+0.1%	-4.2%	
Event Selection	+0.7%	+0.4%	+0.7%	+0.7%	+0.4%	+0.7%
Event Selection	-0.2%	-1.1%	-1.7%	-0.2%	-1.1%	-1.7%
Flavor Tagging	-0.4%	+3.7%	+0.2%	-0.4%	+3.7%	+0.2%
Flavor Tagging	+0.2%	-5.0%	-0.7%	+0.2%	-5.0%	-0.7%
Non uniformity	< 0.1%				< 0.1%	
Combined Systematic Uncertainty	+0.7%	+5.5%	+7.6%	+0.7%	+5.5%	+7.6%
	-0.5%	-6.6%	-7.8%	-0.5%	-6.6%	-7.8%

Table 2. Uncertainties of  $H \to b\bar{b}$ ,  $H \to c\bar{c}$  and  $H \to gg$ 

Analysis with more reliable approaches. Systematic uncertainties considered.

BR 0.57%

2.3%

1.7%

# Study on New Data Sets

# **Outlook of New Data Sets**

- Full simulation events with  $\sqrt{s} = 240$  GeV and new geometry
- Much larger statistics than previous samples: (in 2 fermions and single W semi-lepton sample)
- High performance PID based on LICH
- So far only IIH channel are studied
  - Only consider Higgs and semi-leptonic backgrounds

# Efficiencies

- Efficiency of selecting 2 muons dropped from 87% to 75% (we expect efficiency around 95%)
- Efficiency of selecting 1 electron and positron dropped from65% to 55% (we expect efficiency around 90%)
- Fake leptons in single W semi-lepton sample rising by 50 times(0.1%-5%)
- Wrong PID are used.... Hay be solved soon

# Signal Line Shape of lepton recoil system mass: µµH

Signal Recoil Mass: Described by a Crystal ball function + double sided exponential Head



# Signal Line Shape of lepton recoil system mass: µµH

Signal Recoil Mass: Described by a Crystal ball function + double sided exponential Head



# **Results of Template Fit: µµH**



- Good fit quality
- Slightly higher uncertainty than that in previous study

1.1%

10.5%

5.4%

# Plan of Future

What can we learned from the previous analysis, and in which way we can move forward?

## Lepton Reconstruction with FSR

- Affect on lepton energy/momentum, lepton mass and lepton recoil system mass spectrum
- Affect on the lepton isolation
  - Lower lepton efficiency
- Affect on jet clustering
  - Photons clustered into jets
- We need a robust algorithm to remove radiated photons and recover their energy/momentum into the leptons

# New Strategy on qqH analysis



How might we take this into account?

- Study the qqH->qq+jj by Z and H flavor categories: 9 categories, explore each of them with different mass pair resolution
- Explicitly tag jets with soft leptons(heavy flavor with semi-leptonic decay), try to recovery the missing energy/momentum

## Gluon Jets Issues: Contamination in Heavy Flavor



 We need to not only count the flavors, but also consider their distributions in sub-jet level

## Gluon Issues: H->gg Analysis

Now H->bb/cc/gg are combined as a 'signal' unless the template fit was applied.

- H->bb has the largest fraction, the optimization of event selection might not so good to the other two channels
- Is it possible to apply analysis optimized only for H->gg

### An explicit H->gg analysis



Image is constructed with the energy of all the final state stable particles in an event.



by Wang Yan, Li Zhao and Li Gexing



Each convolutional layer is consisted of 64 or 128 filters with filter size 3\*3 and a ReLU activation.

Each maxpooling layer performed a **2\***2 down-sampling with a stride length of 2.

# Gluon Issues: gluon control sample

- Like other final states, gluon jets control sample to commissioning the distributions of gluon
  - Particle number, flavor etc., need high purity and high efficiency
- Unlike other final states, this is not easy to find
  - Three jets events?

# Systematic Uncertainty of Flavor Tagging

A lot of things need to be measured in data, and compare to MC prediction: Track multiplicity in jets, tracks' impact parameters, secondary vertex variables, B/C-likeness, X<sub>B</sub>-X<sub>C</sub> distributions(correlation of tagging variables)

We don't have data yet.

We need to demonstrate how well we could measure those variables in data.

## A simple example

To Estimate the systematic of b-tagging in  $\mu\mu$  channel

- The uncertainty are directly from the difference in 'data' and MC-prediction of the templates
- Assuming we use ZZ->µµ+bb to calibrate the bb-templates, and MC prediction bias for some systematic uncertainty reason(H->bb should also have such bias).
- Select a ZZ-> $\mu\mu$ +qq sample. The purity of sample is 99.6%, with more than 20k ZZ-> $\mu\mu$ +bb events.
- Estimate how big a difference between MC and data would be detected, with a limit in the data statistic uncertainty and the precision of R<sub>b</sub> prediction
- Estimate the impact to the branch fraction result with such difference

## c-jets control sample

- b-c jets separation is relatively worse than c-light jet separation
- A control sample with low b-jets contamination is required
- Semi-leptonic WW sample might be a choice
- Selection on WW-> μ ν+qq events gives high statistic and high purity control sample

	µ+qq	τ+qq	zz->µµ+qq	Other
Event Yields	2.2M	38k	10k	8.3k
Fraction	97.5%	1.7%	0.43%	0.36%

About half of them contains a c-jets<sub>26</sub>

- tau+qq with tau -> mu decay fakes to WW-> μ ν+qq, but they should have similar W hadronic decay components
- WW semi-leptonic events has a purity over 99%, and

# Summary

- H->bb/cc/gg measurements are very important in understanding quark's mass origin, and are important BM analysis in CEPC
- Previous work with old geometry verified the capability of high precision measurements of H->bb/cc/gg in CEPC
- PID need to be corrected. But the fit seems promising.
- A lot of things to try and improve: leptons reconstruction, jets paring, gluon jets study, flavor tagging calibration...

#### Let's keep up the good work!

# Thank You!

# Backup

### Likelihood Function of 3D Fit

$$\begin{split} & L(M_{\text{recoil}}^{l\bar{l}}, X_B, X_C; \vec{\theta_s}, \vec{\theta_b}, N_{H \to b\bar{b}}^{\text{sig}}, N_{H \to c\bar{c}}^{\text{sig}}, N_{H \to gg}^{\text{bkg}}, N_{\text{irred}\_b\bar{b}}^{\text{bkg}}, N_{\text{irred}\_uds}^{\text{bkg}}, N_{l^+l^-H}^{\text{bkg}}, N_{l^+l^-H}^{\text{bkg}}, N_{redu}) \\ = & P_{\text{sig}}(M_{\text{recoil}}^{l\bar{l}}; \vec{\theta_s})(N_{H \to b\bar{b}}^{\text{sig}} P_{\text{flavor}}^{H \to b\bar{b}}(X_B, X_C) + N_{H \to c\bar{c}}^{\text{sig}} P_{\text{flavor}}^{H \to c\bar{c}}(X_B, X_C) + N_{H \to gg}^{\text{sig}} P_{\text{flavor}}^{H \to gg}(X_B, X_C) + N_{l^+l^-H}^{\text{bkg}} P_{\text{flavor}}^{H \to other}(X_B, X_C)) \\ & + P_{\text{irred}}(M_{\text{recoil}}^{l\bar{l}}; \vec{\theta_b})(N_{\text{irred}\_b\bar{b}}^{\text{bkg}} P_{\text{flavor}}^{\text{irred}\_b\bar{b}}(X_B, X_C) + N_{\text{irred}\_c\bar{c}}^{\text{bkg}} P_{\text{flavor}}^{\text{irred\_c\bar{c}}}(X_B, X_C) + N_{\text{irred}\_uds}^{\text{bkg}} P_{\text{flavor}}^{\text{irred\_uds}}(X_B, X_C))) \\ & + N_{\text{redu}} P_{\text{redu}}(M_{\text{recoil}}^{l\bar{l}}, X_B, X_C). \end{split}$$

# **3D fit Plots**



# ToyMC Test of 3D Fit



# H->gg conventional results



# Correlation of FT on two jets

- We need to tag two jets in a events: the tagging/mis-tagging rate will be correlated
- Correlation is small (up to a few percent), but need to be considered to achieve high precision

#### The correlation of FT are considered in R<sub>b</sub> measurement in ALEPH:

- Phys.Lett.B 401(1997) 163
- Phys.Lett.B 401(1997) 150

#### Think about di-jet sample:

- $f_s = R_b \boldsymbol{\xi}_b + R_c \boldsymbol{\xi}_c + (1 R_b R_c) \boldsymbol{\xi}_{uds}$
- $f_d = R_b \mathcal{E}_b^2 (1 + q_b) + R_c \mathcal{E}_c^2 + (1 R_b R_c) \mathcal{E}_{uds^2} + \mathcal{E}_x$ : Efficiency for flavor x
- f<sub>s</sub>, f<sub>d</sub> measured from data
- R<sub>c</sub> from other measurement
- $\mathbf{E}_{c}$ ,  $\mathbf{E}_{uds}$  and  $\mathbf{g}_{b}$  taken from MC
- Work out R<sub>b</sub> and **E**b

- **f**<sub>s</sub> : Single Hemisphere Tagged
- **f**<sub>d</sub> : Both Hemisphere tagged
- **g**<sub>b</sub>: Hemispheres Correlation

 $Q_b < \xi_b >^2 = < \xi_b^2 > - < \xi_b >^2$ 

In less symmetric case:  $Q_b < \varepsilon_1 > < \varepsilon_2 > = < \varepsilon_1 \varepsilon_2 > - < \varepsilon_1 > < \varepsilon_2 >$