# FireWork evaluation from a forecaster point of view

**Bruce Ainslie & Rita So**
**Environment & Climate Change Canada**

NW-AIRQUEST 2020 Annual Meeting
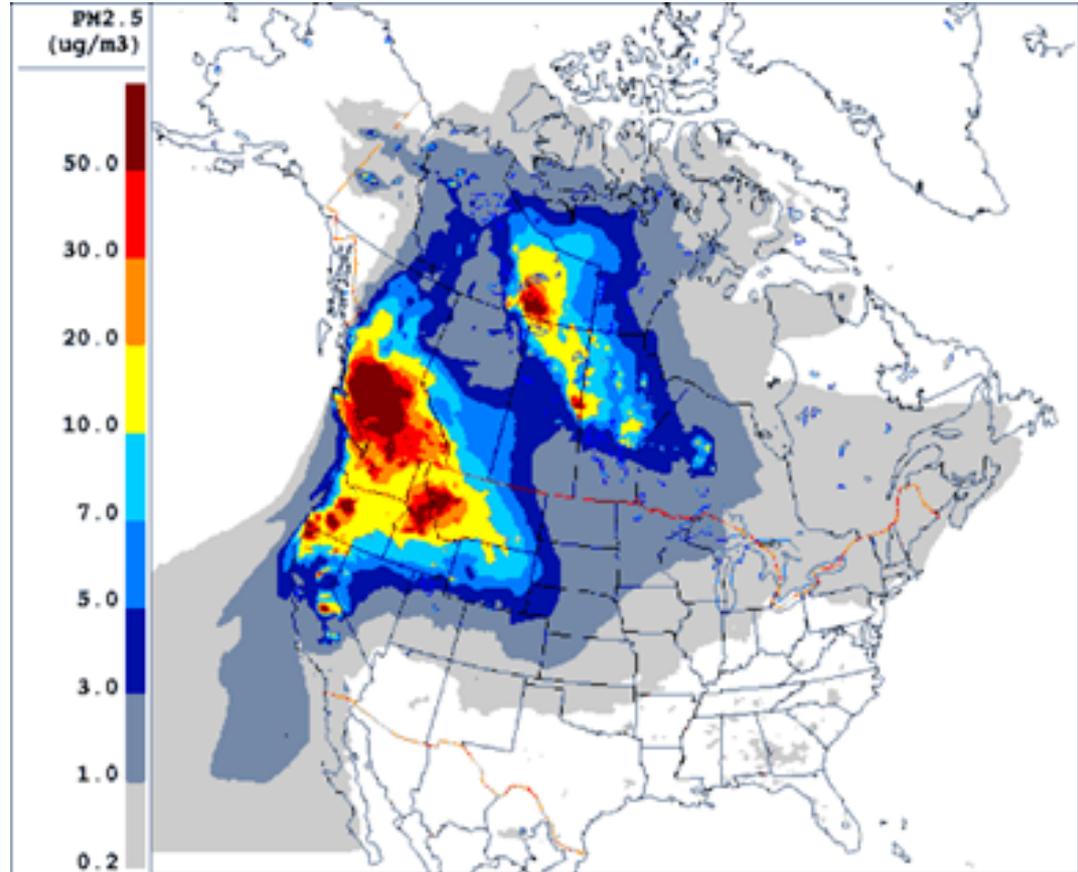
June 10, 2020

NASA Aqua MODIS imagery on Aug 1, 2017.
Source: LANCE/EOSDIS MODIS Rapid Response Team

# ECCC's FireWork System

## https://weather.gc.ca/firework

- Twice daily (00z/12z) during fire season (Apr- Oct)
- Incorporated into the ECCC Regional Air Quality Deterministic Prediction System
- NRT fire data from Canadian Wildland Fire Information System (based on NOAA/NASA satellite info.)
- Products:
  - PM2.5 & PM10 maps
  - AQ point forecasts
  - 24-hr accumulated PM2.5

# How to evaluate performance?
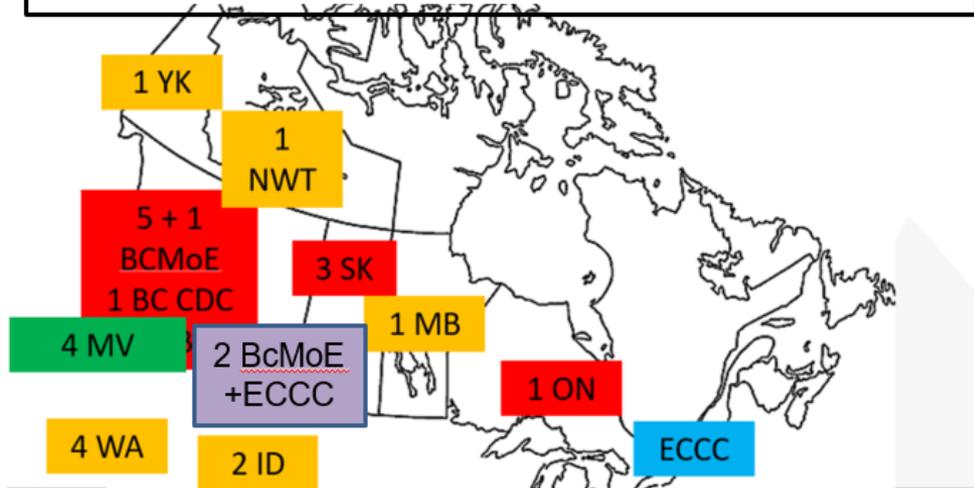# How do we make use of forecast guidance?



I. Informal meeting with MV
II. https://www.surveymonkey.com/r/ZKXHJG8 (original)
III. Teleconference with ECCC
IV. https://www.surveymonkey.com/r/FL7KKT9 (updated)
V. Informal face to face meetings with forecasters

1 YK
1 NWT
5 + 1 BCMoE 1 BC CDC
3 SK
4 MV
2 BcMoE +ECCC
1 MB
1 ON
4 WA
2 ID
ECCC

**FireWork User Survey**

# FireWork User Survey - Takeaways

## Survey shows emphatically **model guidance is useful**

### Four key **perspectives**:

I. Jurisdiction-based, but never at a monitor
II. Event-based: **25 ug/m3**
III. Missing an event >> False alarm
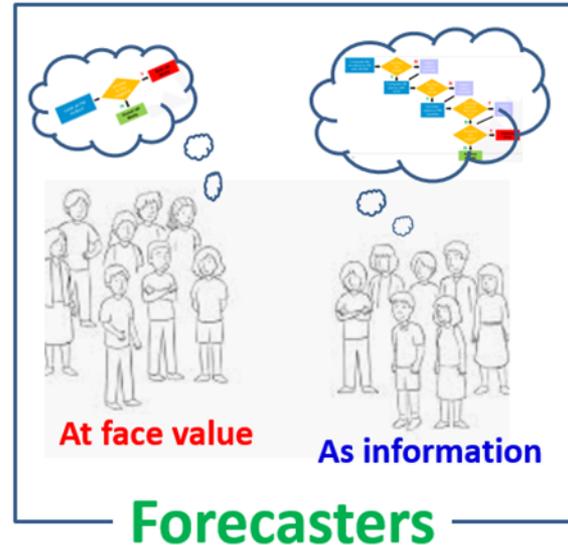IV. 1st day forecast (0 – 24 hrs) is the MOST IMPORTANT



### Two **types** of users:

I. Take forecast "at face value"
II. Take forecast "as information"



At face value

As information

**Forecasters**

# Model Evaluation Framework

Use survey results to design evaluation framework and test ECCC's old (**FEPS**) and new (**CFFEPS**) operational wild fire smoke models over 2016-2018 fires seasons
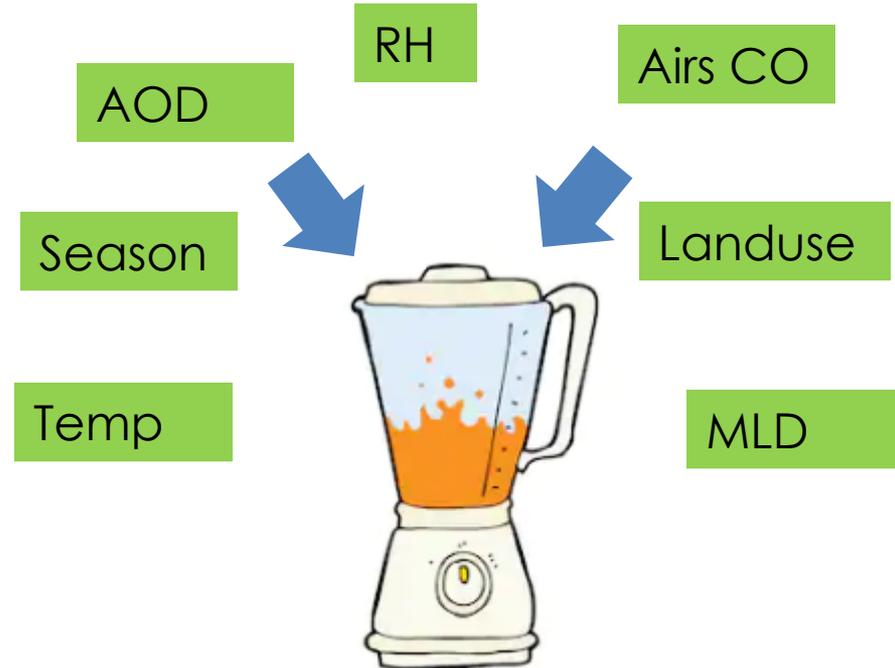


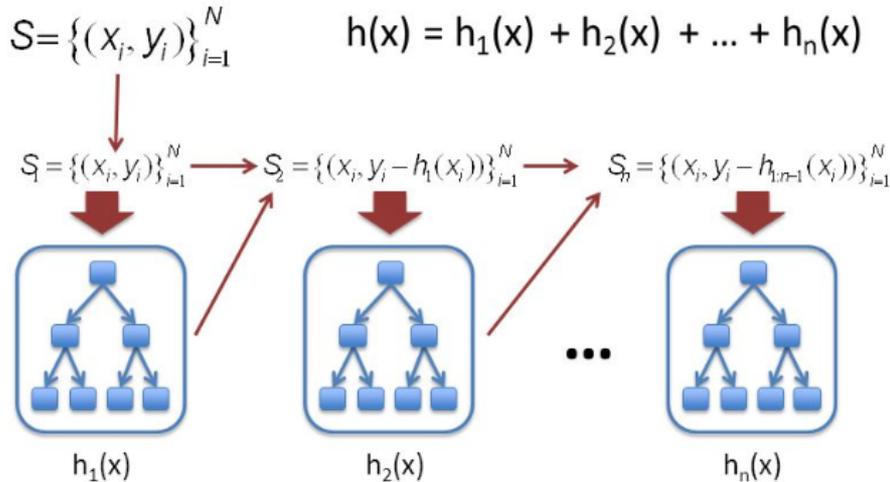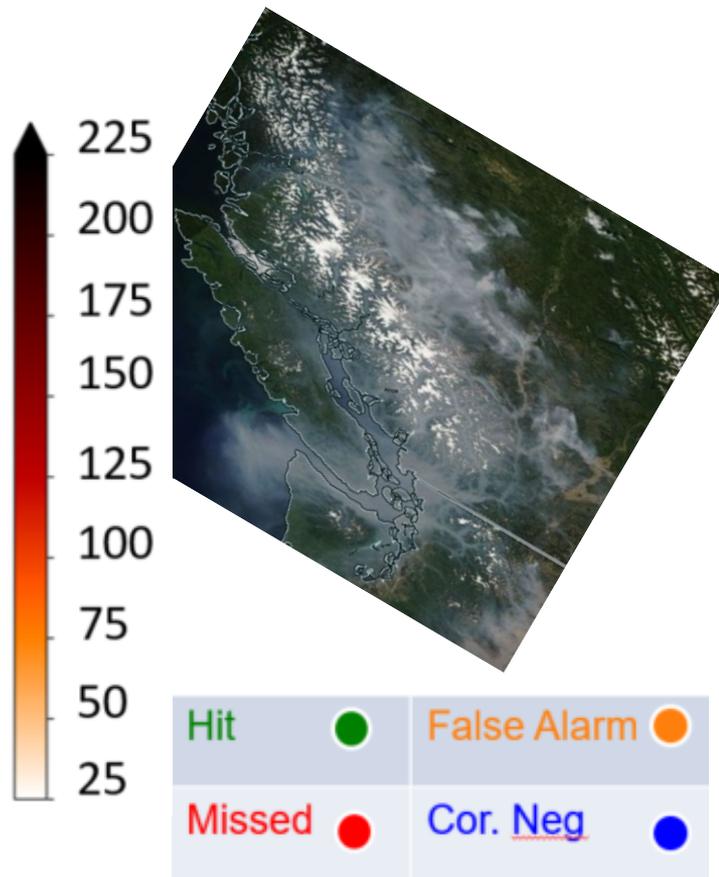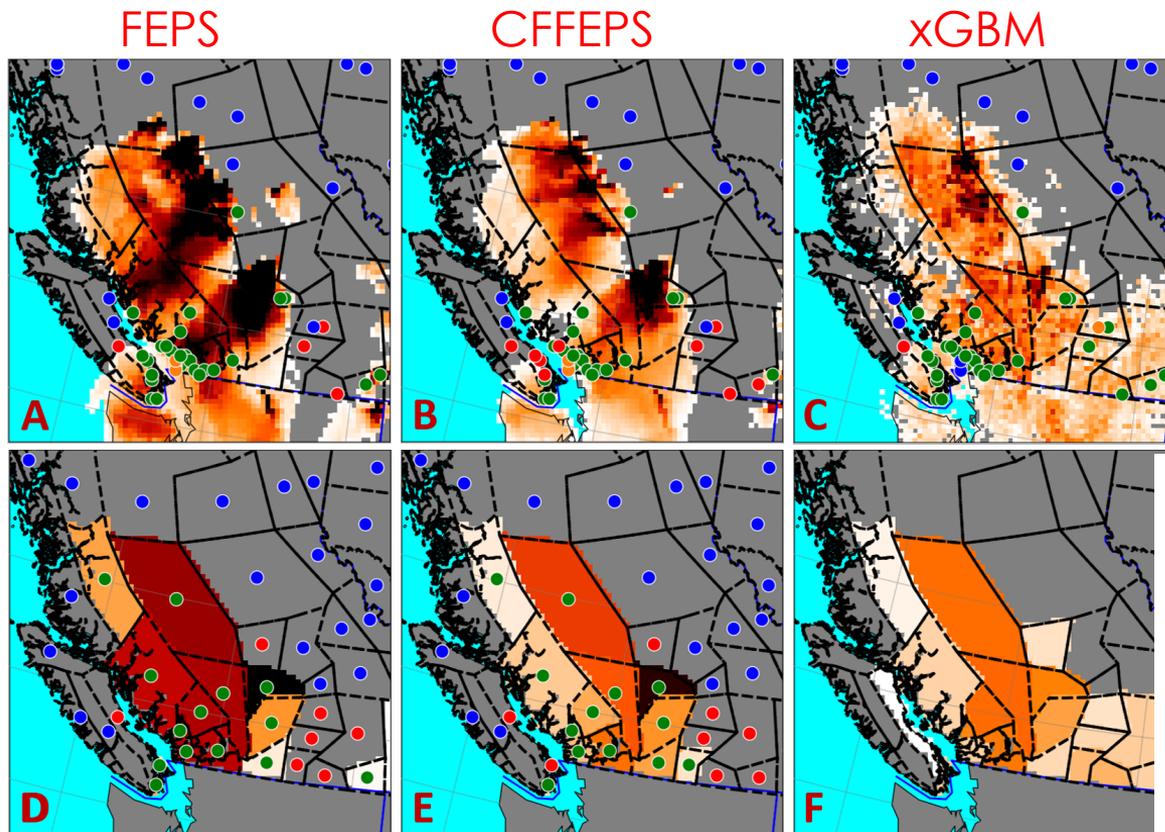Meteorology (GEM) + Fire Emissions → FireWork

**Gridded Spatial Observation dataset**
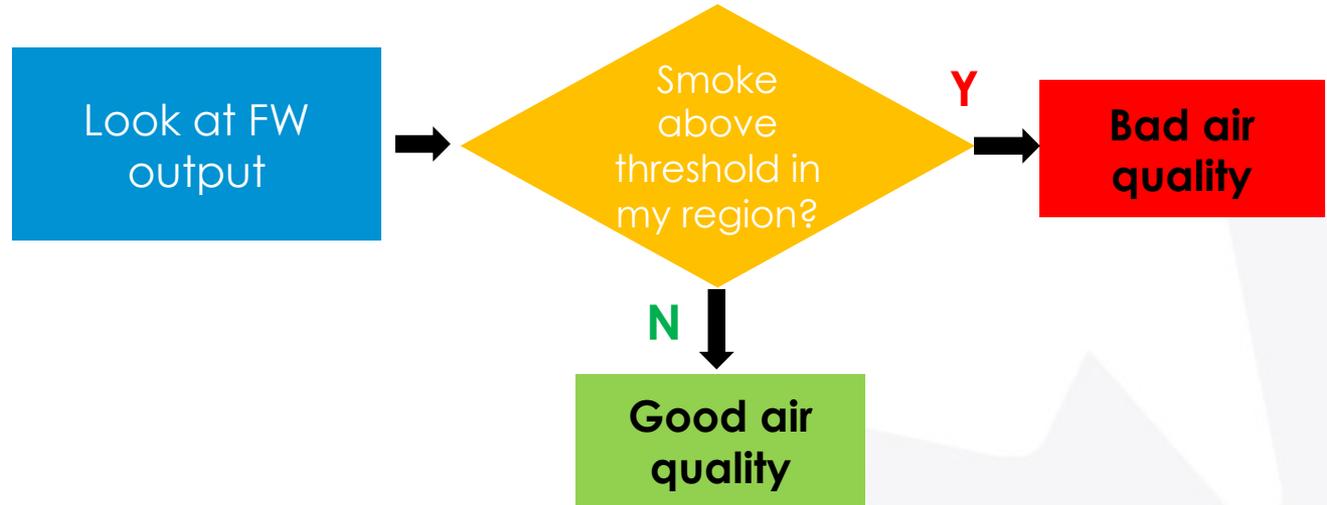
# Detour: Extreme Gradient Boosting (xGBM)

- Use machine learning to downscale satellite data any many other dataset to construct daily 24-hr PM2.5
- GBM predicts in the form of an ensemble of simple models
- It builds model in a stage-wise fashion

$$S = \{(x_i, y_i)\}_{i=1}^{N} \qquad h(x) = h_1(x) + h_2(x) + \dots + h_n(x)$$

$$S_1 = \{(x_i, y_i)\}_{i=1}^{N} \rightarrow S_2 = \{(x_i, y_i - h_1(x_i))\}_{i=1}^{N} \rightarrow S_n = \{(x_i, y_i - h_{1:n-1}(x_i))\}_{i=1}^{N}$$

$h_1(x)$    $h_2(x)$    $\dots$    $h_n(x)$

RH

Airs CO

AOD

Season

Landuse

Temp

MLD

# Detour: Spatial Footprints August 2nd 2017

# Model Evaluation I:
## "At face value" decision making process



Information

Evaluation

Update

Look at FW output

Smoke above threshold in my region?

Y

Bad air quality

N

Good air quality

# At face value event-based verification: contingency table

Event-based + at face value → use standard 2x2 contingency metrics

Would like metric(s) to be:

1. Independent of **forecast situation** (e.g. separate *forecast system* from *forecast situation*)

2. Punish MISS more than FA

→ Unfortunately, treated misses differently than false alarms is **value** not a **quality**

# Peirce Skill Score

- Large number of ways of scoring contingency table

- Table can be expressed in terms of 3 parameters:

  1. **H: hit rate a/(a+c)**
  2. **F: false alarm rate b/(b+d)**
  3. **s: base rate (a+c)/N**

We have chosen **PSS** = **H** − **F** since:
  1. independent of base rate
  2. Provides a link between forecast quality and value

|  | Event Observed | | Totals |
|---|---|---|---|
| Event Forecast | a (hits) | b (false alarms) | a+b |
|  | c (misses) | d (correct neg.) | c+d |
| Totals | a+c | b+d | N |

# PSS: CFFEPS vs FEPS
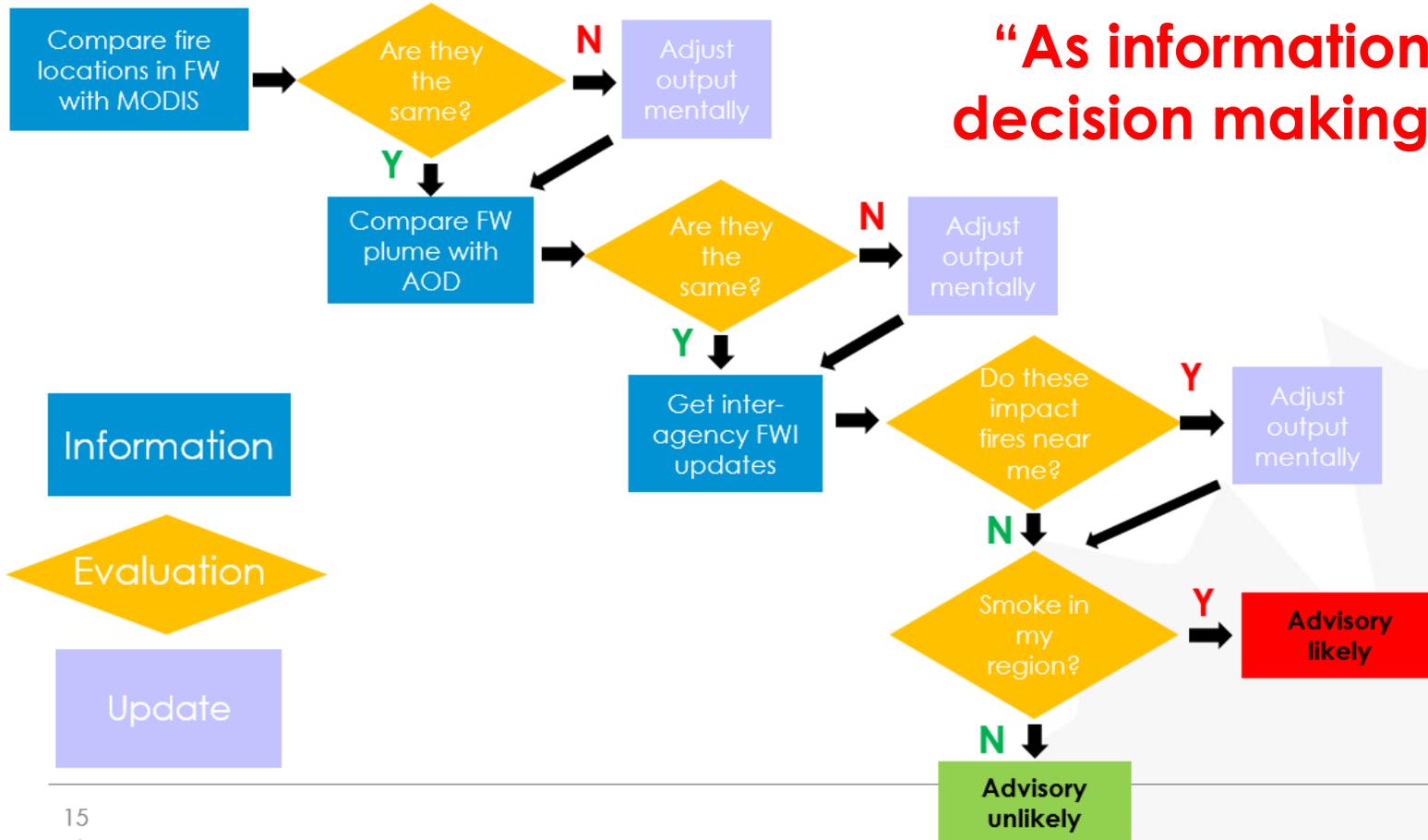## by MSC forecast zone

Wilcoxon signed-rank test suggests the two estimates for MSC Zone PSS scores do not come from the same population (p<0.01).
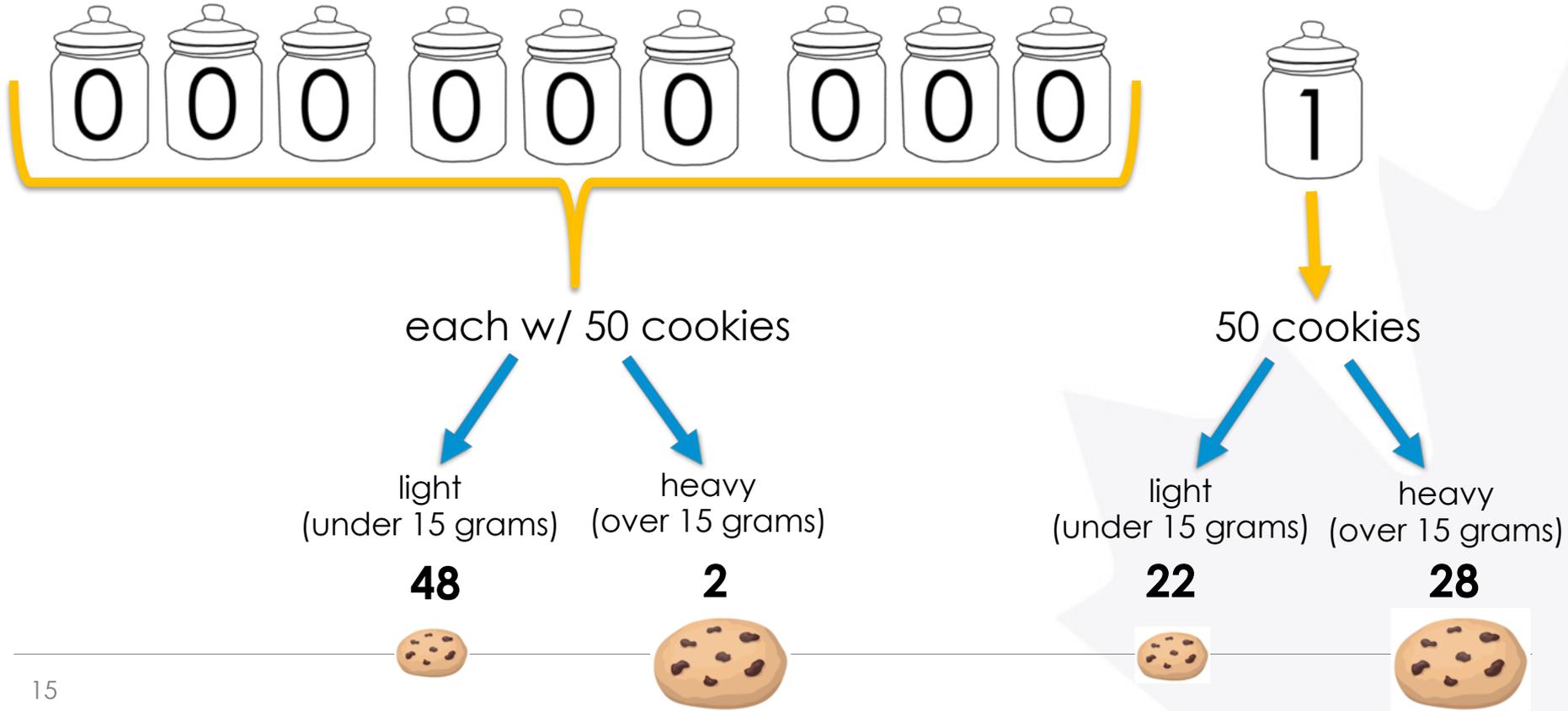➢ Feps provided more 'value' than Cffeps



PSS

# Model Evaluation II: "As information source" decision making process



| | |
|---|---|
| Compare fire locations in FW with MODIS | → |

Are they the same? — N → Adjust output mentally

Y ↓

Compare FW plume with AOD → Are they the same? — N → Adjust output mentally

Y ↓

Get inter-agency FWI updates → Do these impact fires near me? — Y → Adjust output mentally

N ↓

Smoke in my region? — Y → **Advisory likely**

N ↓

**Advisory unlikely**

**Information**

**Evaluation**

**Update**

**Bayesian Refresher**

# Detour: Which cookie jar?



each w/ 50 cookies

50 cookies

light
(under 15 grams)

heavy
(over 15 grams)

light
(under 15 grams)

heavy
(over 15 grams)

48

2

22

28
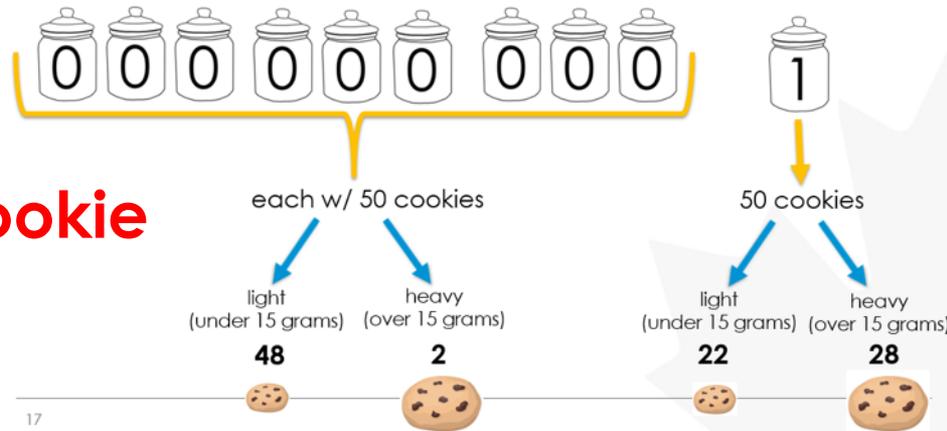
A single cookie is taken from one of the ten jars.

**What is the probability the cookie was drawn from jar 1?**

**p(J=1) = 0.1**

each w/ 50 cookies
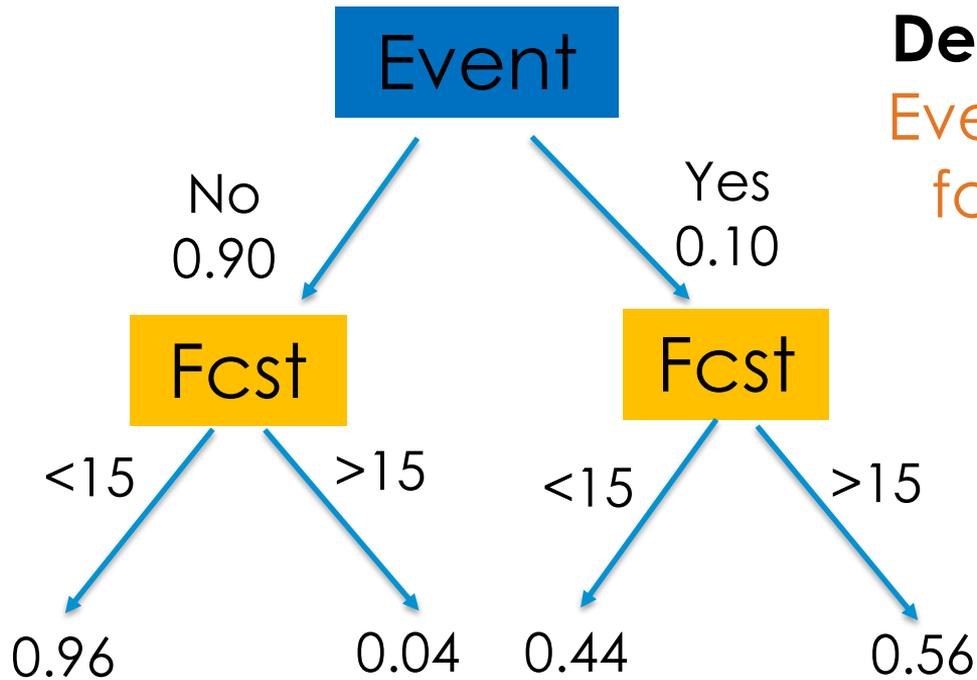
50 cookies

light
(under 15 grams)

heavy
(over 15 grams)

light
(under 15 grams)

heavy
(over 15 grams)

**48**

**2**

**22**

**28**

17

Now, if I tell you the drawn cookie weighted over 15 grams – **how does this information help refine your guess?**

**p(J=1|heavy) =** p(heavy|J=1)/p(heavy)   x   p(J=1)

**=** (28/50)/(0.1*28/50 + 0.9*2/50)  x  0.1

**0.61 = 6.1 x 0.1**

| Updated guess | New information | Initial guess |

**Detour:** Chance of an Event today given the forecast PM2.5 was greater than 15 ug/m3?

Event

No 0.90          Yes 0.10

Fcst                    Fcst

<15        >15        <15        >15

0.96       0.04       0.44       0.56

p(**Event**|**Fcst>15**) = p(Fcst>15|Event)/p(Fcst>15)x BaseRate
= 6.1 x 0.1
= **0.61**

17

# As Information source: convert forecast → p(Event)

**Binomial (Logistic) Regression**

$$P(Event|fcst) \sim \frac{1}{1 + exp\{-1[bo + b1 * conc + b2 * \Theta + b3 * met]\}}$$

**Our regression using GAM**

$$\log\left(\frac{P}{1-P}\right) \sim f(conc) + g(\Theta) + j(NCEP)$$

Event probabilities as a function of forecast smokiness

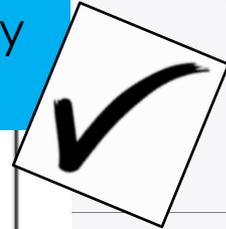**CFFEPS Logistic Regression**
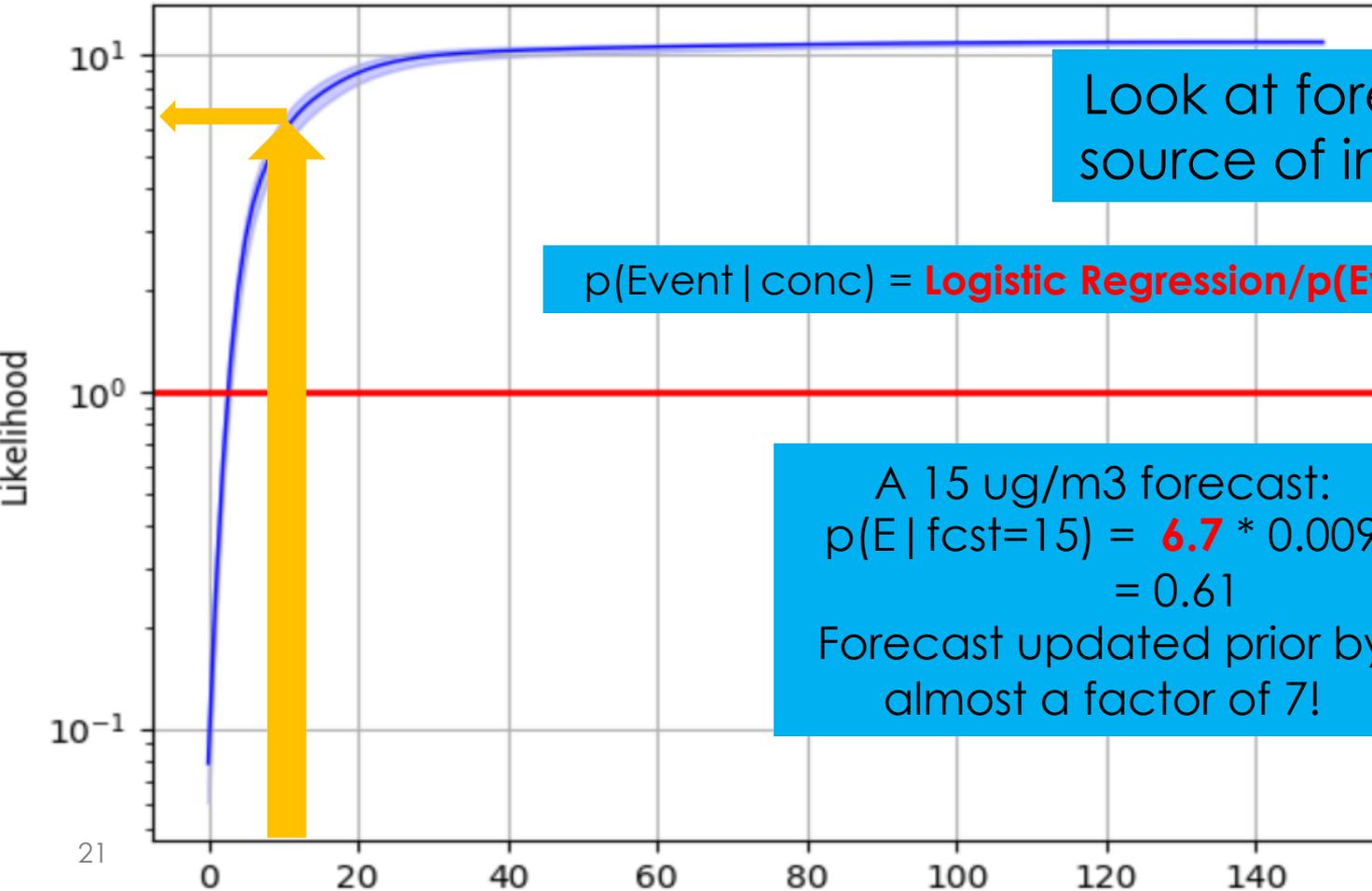(showing regression for
$\theta = 0.009$ base rate only)

19

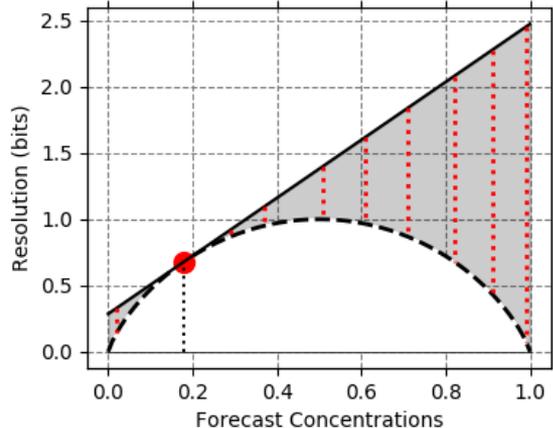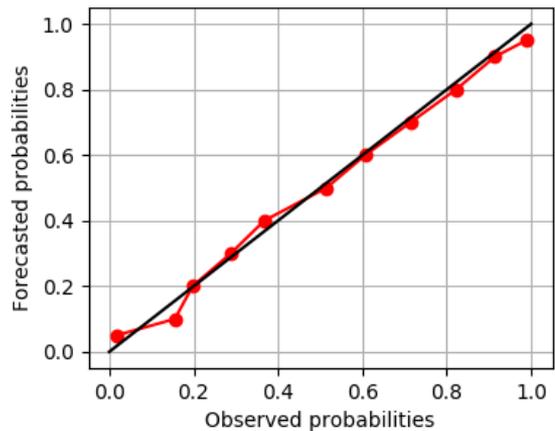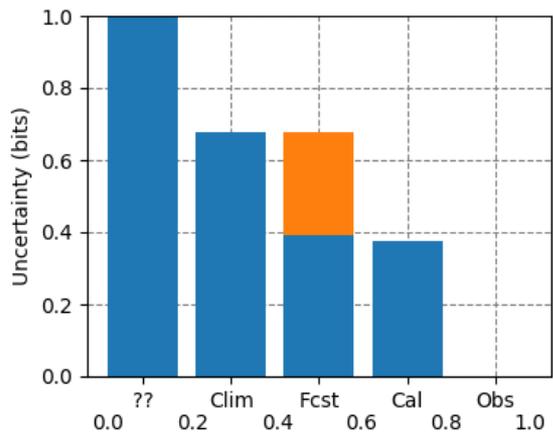Event probabilities as a function of forecast smokiness

For a 15 ug/m3 forecast:
p(E|fcst=15) = 0.61

**There is a 61% chance of a missed event !**
(when taken at face value)

Mean likelihood

Look at forecast as a source of information

p(Event|conc) = **Logistic Regression/p(Event)** * p(Event)

A 15 ug/m3 forecast:
p(E|fcst=15) = **6.7** * 0.009
= 0.61
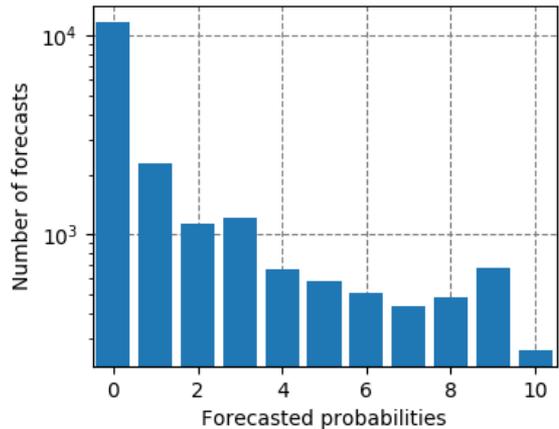Forecast updated prior by almost a factor of 7!

## Question:
## *How to tell which model supplies more information?*

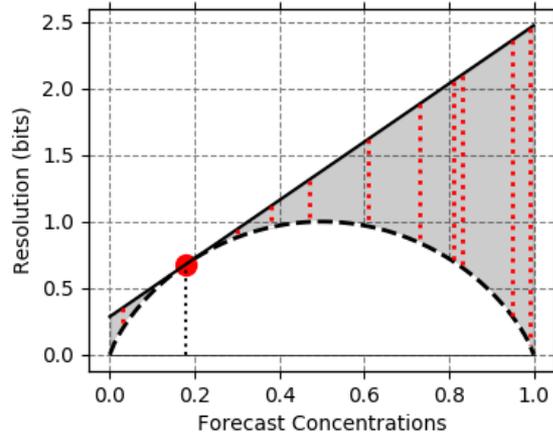1. Ranked Mutual Information Skill Score (RMIS)
2. Divergence Skill Score (DSS)
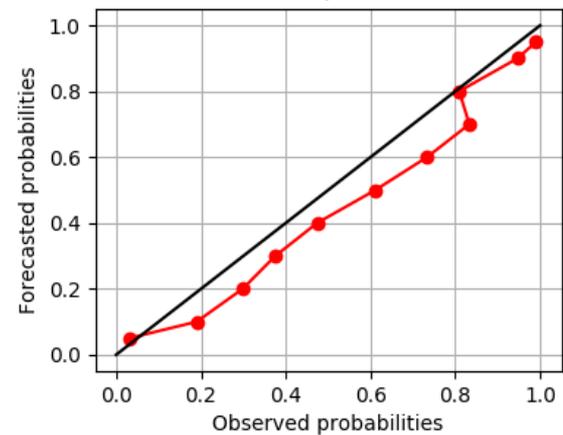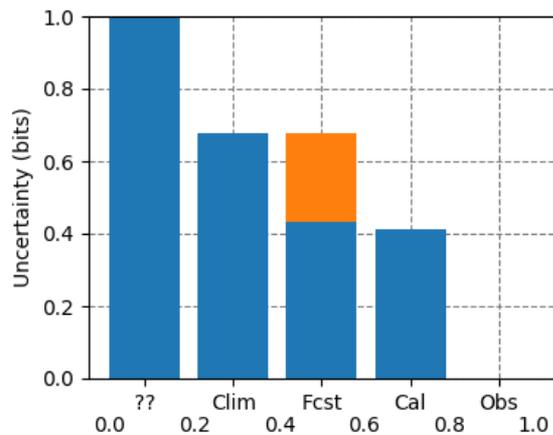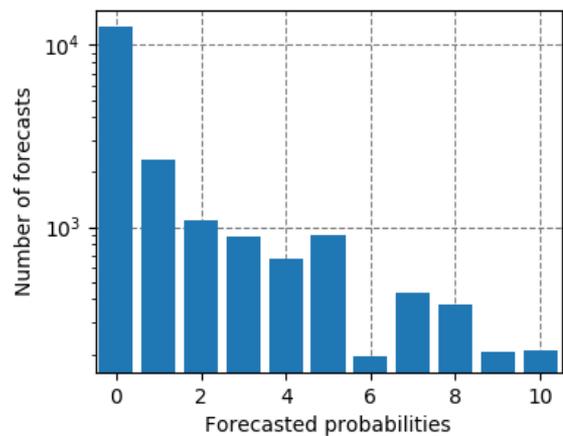3. Ignorance Skill Score (IGN)

Measures information in forecast conveyed to user
→ how much more you know about the future after seeing the forecast, but starting from base rate
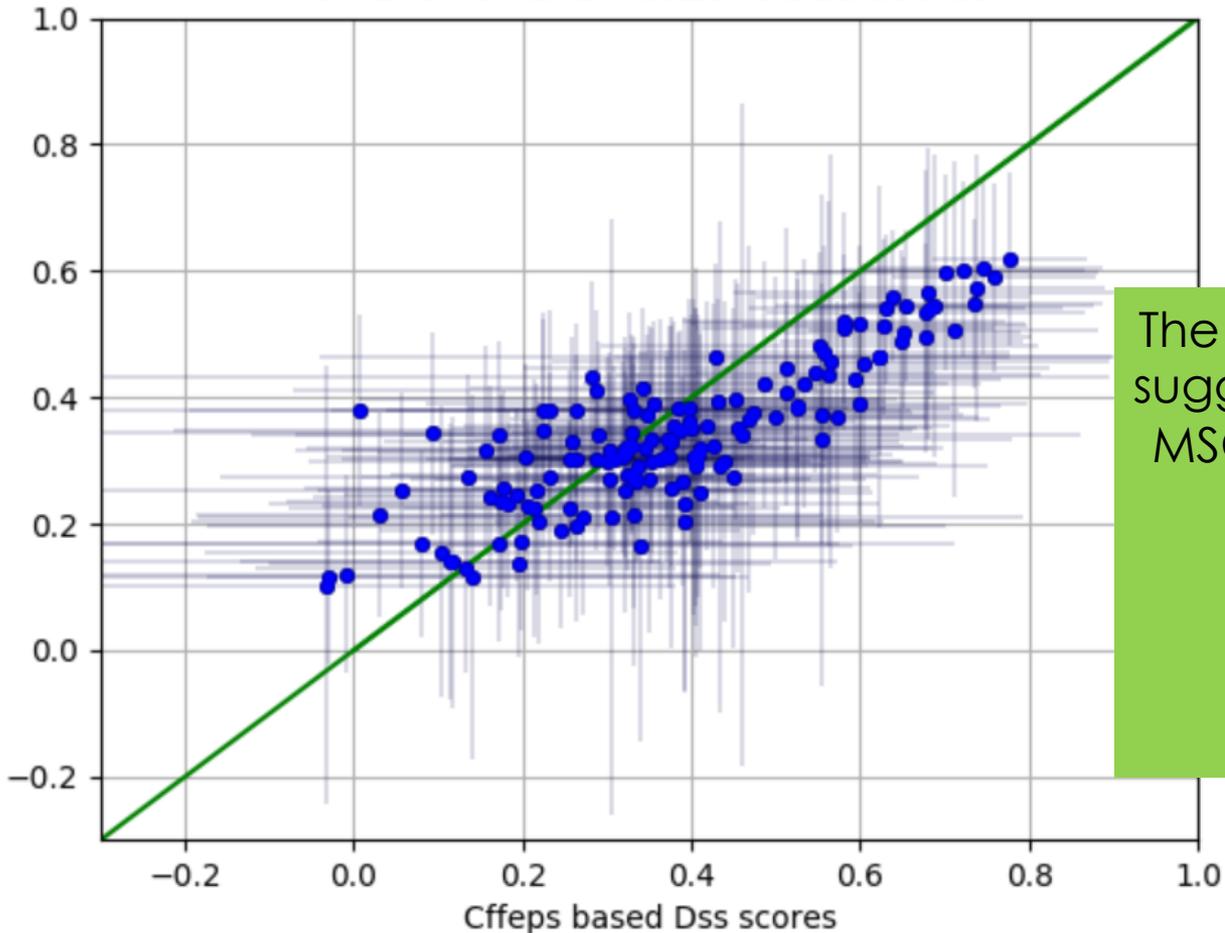
DSS:0.42 Ds:0.39 Rel:0.02 Res:0.30 Unc:0.68



**CFFEPS**

DSS:0.36 Ds:0.43 Rel:0.02 Res:0.26 Unc:0.68

**FEPS**

CFFEPS vs FEPS MscZone DSS scores

The Wilcoxon signed-rank test suggests the two estimates for MSC Zone DSS scores do not come from the same population (p<0.01).
 ➢ CFFEPS provides more information to forecaster than FEPS

# Conclusions

- Guided by survey: 2 interpretations of model guidance, both focussed on events

- Evaluation at the forecast zone-scale and over all of Western Canada via xGBM

- "At face value" analysis via PSS score: FEPS a little more skillful/valuable than CFFEPS

- "As information" analysis via DSS: CFFEPS a little more informative than FEPS

*Did we miss something in how guidance is used?*