# Weakly supervised methods for LHC analyses

Thorben Finke (finke@physik.rwth-aachen.de)
RWTH Aachen University
Institute for Theoretical Particle Physics and Cosmology

# What is the strength of weakly supervised methods?
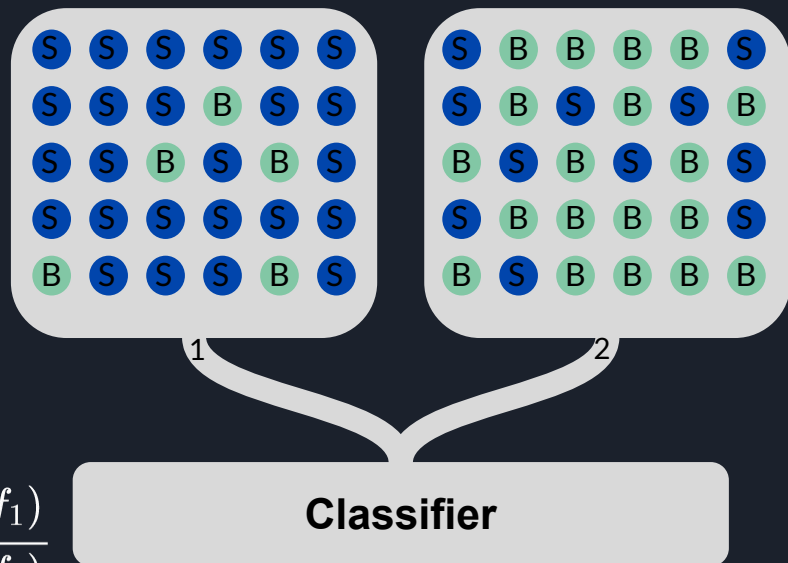
**No need for truth level labels**

$\Rightarrow$ Can train directly on data
  - ○ Avoid systematic uncertainties arising when applying a NN trained on Monte Carlo to experimental data
$\Rightarrow$ Signal (and background) model independence

# Classification without labels (CWoLa)

- Two samples **M₁ and M₂** with signal fractions $f_1$ and $f_2$ **with $f_1 > f_2$**
- Same background and signal distributions in $M_1$ and $M_2$
- ⇒ Optimal classifier for $M_1$ and $M_2$ also optimal for signal (S) and background (B)
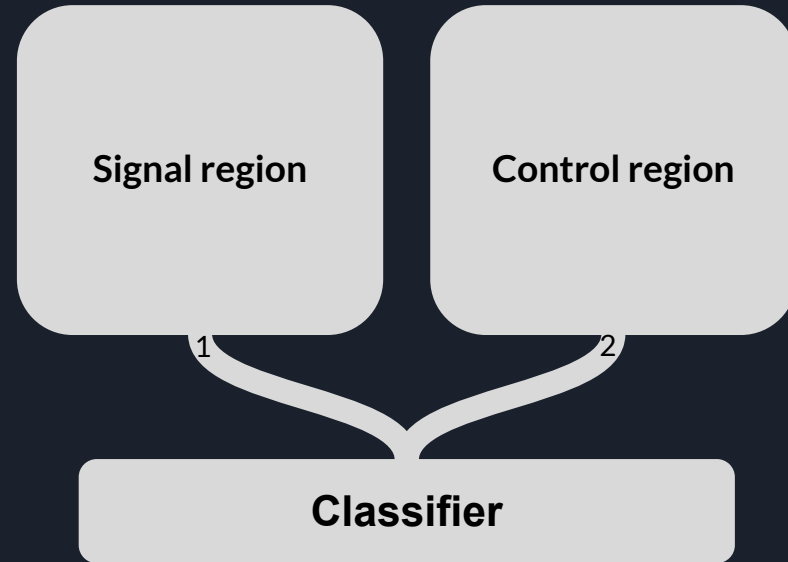


**Classifier**

$$L_{M_1/M_2} = \frac{p_{M_1}}{p_{M_2}} = \frac{f_1 p_S + (1 - f_1)p_B}{f_2 p_S + (1 - f_2)p_B} = \frac{f_1 L_{S/B} + (1 - f_1)}{f_2 L_{S/B} + (1 - f_2)}$$

$$\partial_{L_{S/B}} L_{M_1/M_2} = \frac{(f_1 - f_2)}{\left(f_2 L_{S/B} + 1 - f_2\right)^2} > 0$$
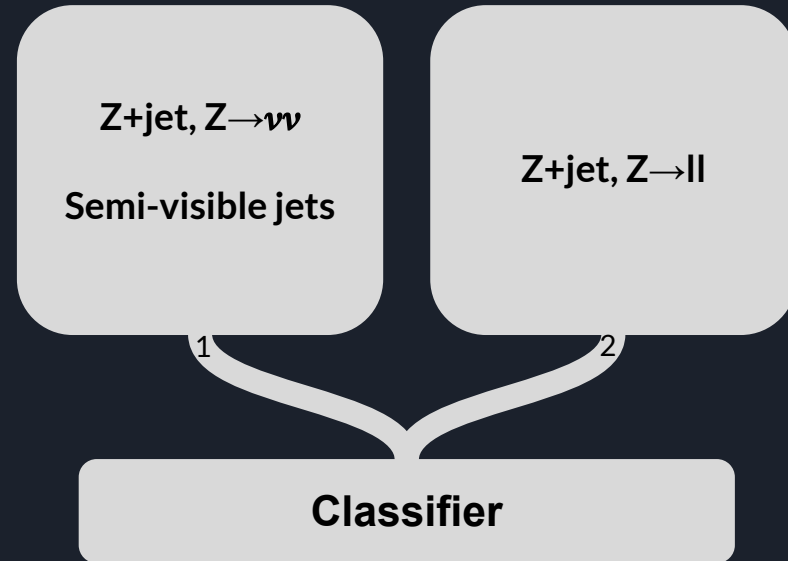
# How to use this for a physics analysis?

- Use control and signal regions as $M_1$ and $M_2$
- Need to ensure same distribution of features ($x$) for background and signal in the two regions
- CWoLa gives sensitivity to differences in the two regions
- Use control region as proxy for the classifier behavior on background in the signal region
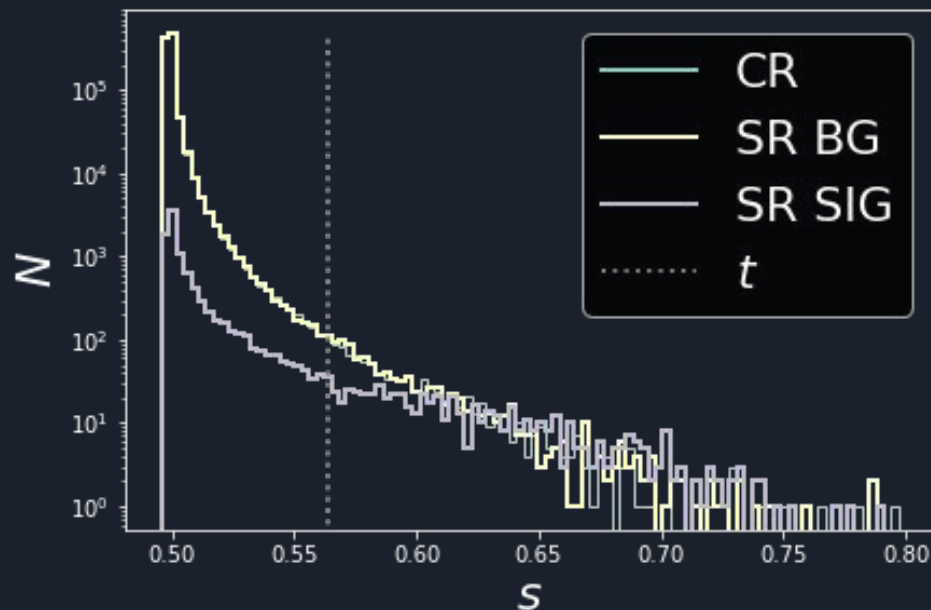
Signal region

Control region

1

2

**Classifier**

# Example: mono-jet search for finding semi-visible jets [arXiv2204.11889](#)

- Signal region with energetic jet recoiling against missing energy
- Semi-visible jets from strongly interacting dark sector as example
  - One jet stays invisible
- $O(10^6)$ background events
- Classify according to jet properties

Z+jet, Z→$\nu\nu$

**Semi-visible jets**

Z+jet, Z→ll

1

2

**Classifier**

# Classifier output

- Peak at ~0.5
  - Expected from indistinguishable background
- Background in signal and control region follow same distribution
- Choose a threshold based on control region
  - Set to keep 0.1 % (1000 events)
- Beyond threshold significant enhancement of S/B
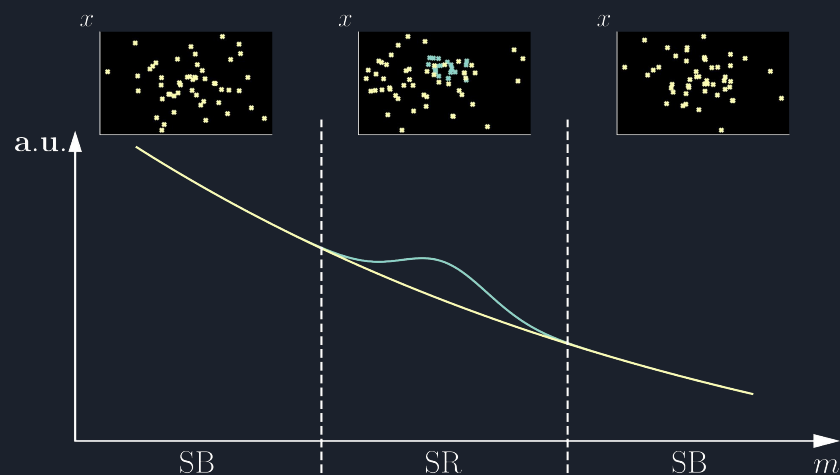
# Results using only main background (Z+jet)

- CWoLa does not introduce fake signal
- High sensitivity beyond current ATLAS limits (<40k events at 95 % CL)

| $f_1$ | $n^{SR}$ | $n^{SIG}$ | stat. sign. |
|-------|----------|-----------|-------------|
| 0 % | 1048 | 0 | 1.07 |
| 0.6 % | 1306 | 247 | 6.84 |
| 1 % | 1666 | 625 | 14.89 |

# CWoLa for bump hunts

Use sidebands (SB) around resonance to estimate background estimation of auxiliar features $x$

1. Use SB data as CR/M$_1$
    a. $x$ must be uncorrelated with $m$
2. Interpolate background features from SB into the SR and use that estimate as CR, e.g. via conditional density estimation -> CATHODE



Recreated from arXiv2109.00546

# Example LHCO2020 R&D data set

LHCO2020

Benchmark data set for anomaly detection

- W -> XY and X/Y -> qq
- $m_W$=3.5 TeV, $m_X$=0.5 TeV, $m_Y$=0.1 TeV
- $m_{jj}$ as resonant feature
- Auxiliary features for the classifier
  $m_{j1}$, $m_{j2}$ - $m_{j1}$, $\tau_{21}^{j1}$, $\tau_{21}^{j2}$
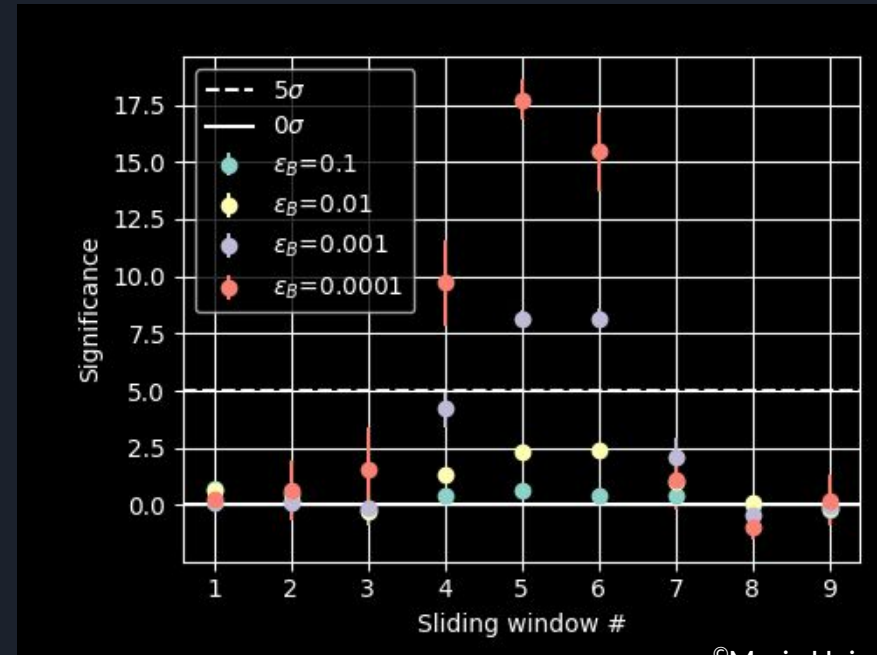
For unknown resonant mass:
    divide into several regions and repeat

# Results of a CATHODE like scan through $m_{jj}$

- Stronger cuts (smaller $\varepsilon_B$) yield higher significance
- Weaker cuts suffer from systematic uncertainty
- Starting from $2.2\,\sigma$, we reach a significance improvement of ~10
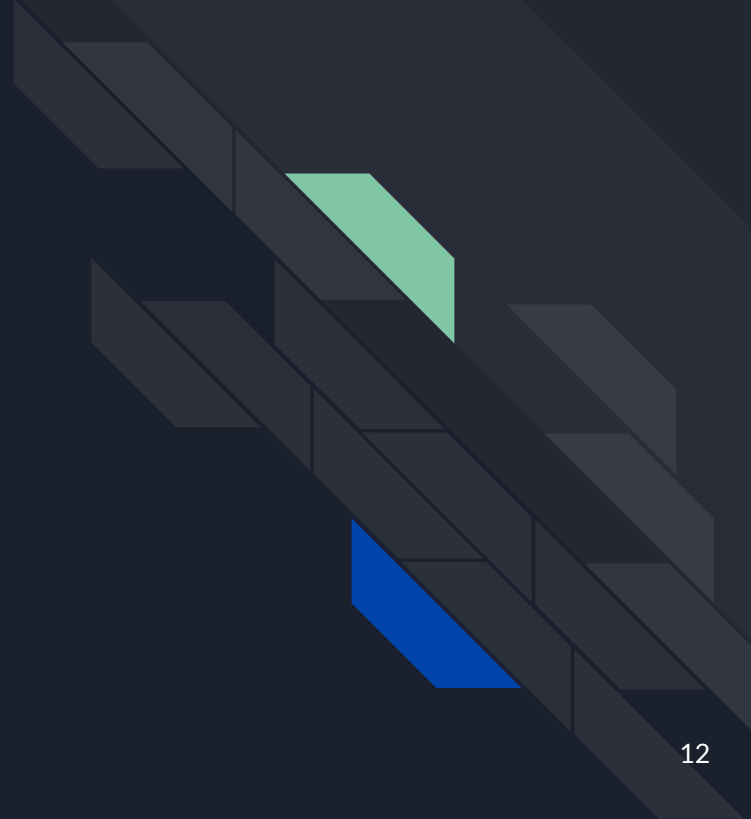


©Marie Hein

# Conclusion

- Weakly supervised methods need no truth level labels
    - ⇒ Avoid systematics from differences in Monte Carlo and data
    - ⇒ Can be applied directly on data
    - ⇒ Are model agnostic and sensitive to a variety of signals
- CWoLa is sensitive to any difference in control and signal region
    - ⇒ Can be used to check validity of the control region
- Setting limits only possible for benchmarks, no model independent limits!

# Backup

# The ATLAS mono-jet search [arXiv2102.10874](arXiv2102.10874)

**Selection cuts:**

- $E_T^{miss} > 200$ **GeV**
- leading AK4 jet with $p_T > 150$ GeV and $|\eta| < 2.4$
- < 4 additional jets with $p_T > 30$ GeV and $|\eta| < 2.8$
- $\Delta\phi(p_T^{jet}, E_T^{miss}) > 0.4$
- lepton veto

**SM backgrounds:**

- **Z+jet** production with invisibly decaying Z (61 %)
- W+jet production with leptonically decaying W and non-identification of the charged lepton (31 %)
- Top quark production (3.5 %)
- Di-boson production (2 %)

Resulting in **$O(10^6)$ background events** and a model agnostic limit of 40k additional events at 95 % CL

# Results using also additional backgrounds

| $r_{tt}^{CR}$ | $r_{VV}^{CR}$ | $n^{SR}$ | $n^{DM}$ |
|---|---|---|---|
| 0 % | 0 % | 4383 | 223 |
| 2.8 % | 1.6 % | 1465 | 456 |
| 3.5 % | 2.0 % | 1686 | 633 |

- Added 3.5 % top and 2 % di-boson background to 1 % signal in signal region
- Ignoring additional backgrounds in control region leads to wrong signal
- Matching the background perfectly recovers the previous performance
- Not matching the background perfectly decreases performance, but does not spoil it completely ⇒ Control region does not need to be perfect

# Results of scan through $m_{jj}$ with SB as CR

- Larger systematic uncertainties reduce the sensitivity compared to CATHODE
- Hardly scratching 5 sigma in wrong window