# Cuckoo!

## Automated Outlier Analysis for Revit

TT Hackathon - 2021

Our Team!

Alex Ogata

Sarvesh Sp

Jim Peraino

Zhenxiang Huang

Matthew Breau

# Problem definition

Uncaught errors in models cost real $$

**Q:** How do you identify when some feature is likely to be an error?

(without hard-coding a huge set of rules)

# What is an outlier?

In statistics, an **outlier** is a data point that differs significantly from other observations.[1][2] An outlier may be due to variability in the measurement or it may indicate experimental error; the latter are sometimes excluded from the data set.[3] An outlier can cause serious problems in statistical analyses.

# What is an outlier?

# Outliers can be hard to spot

**Cell B31:** Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 @ HA_1249_550TenthAVE_Arch_detached 211016_MbreauCNSZQ.rvt

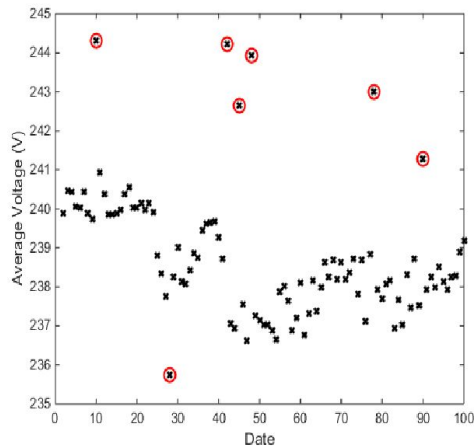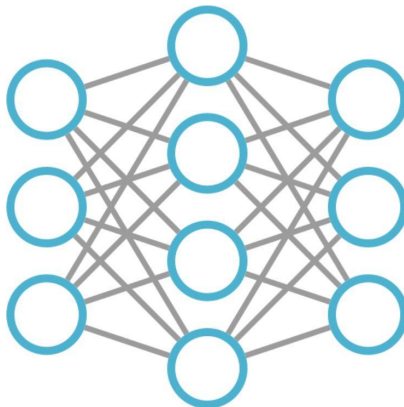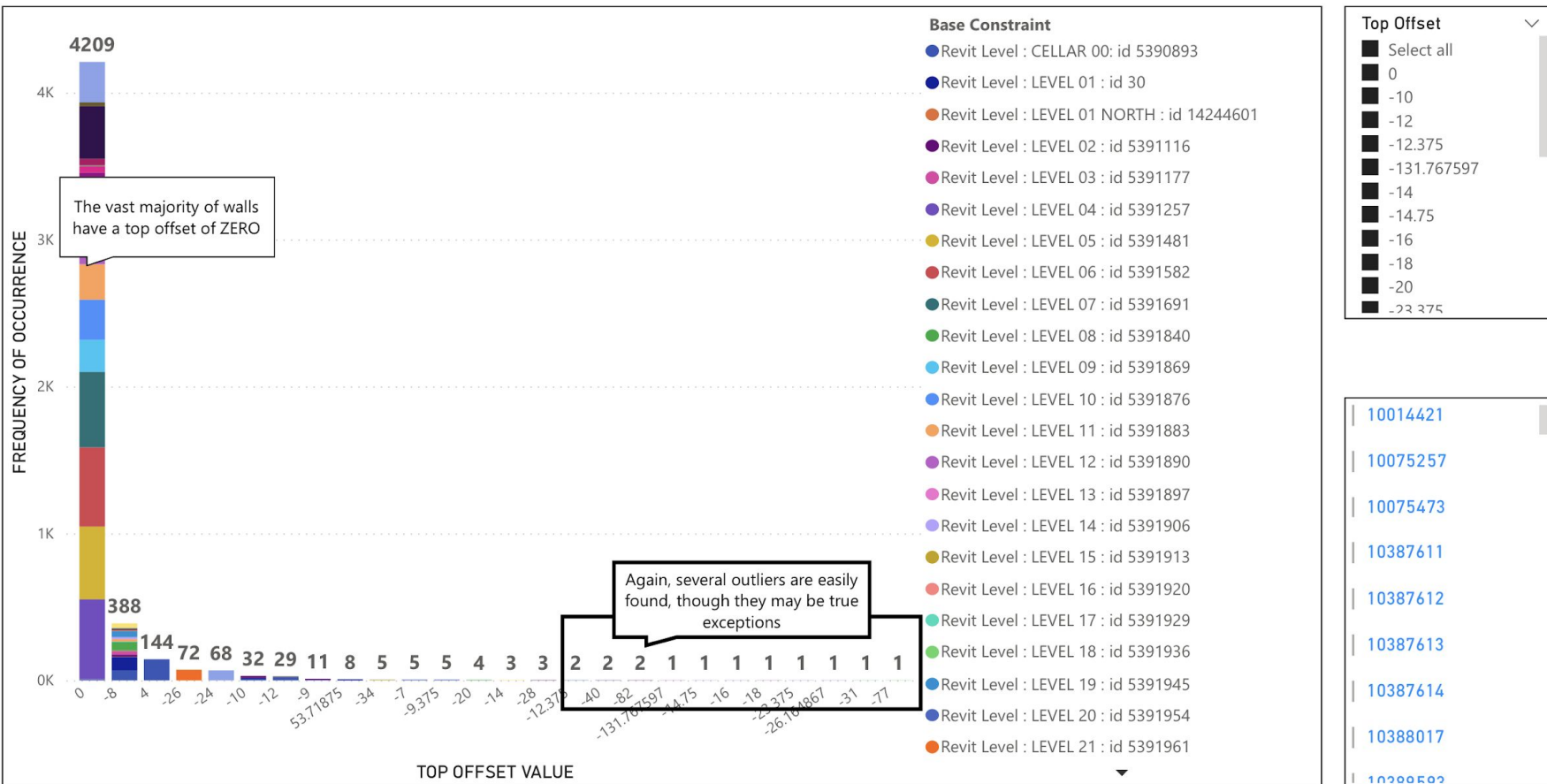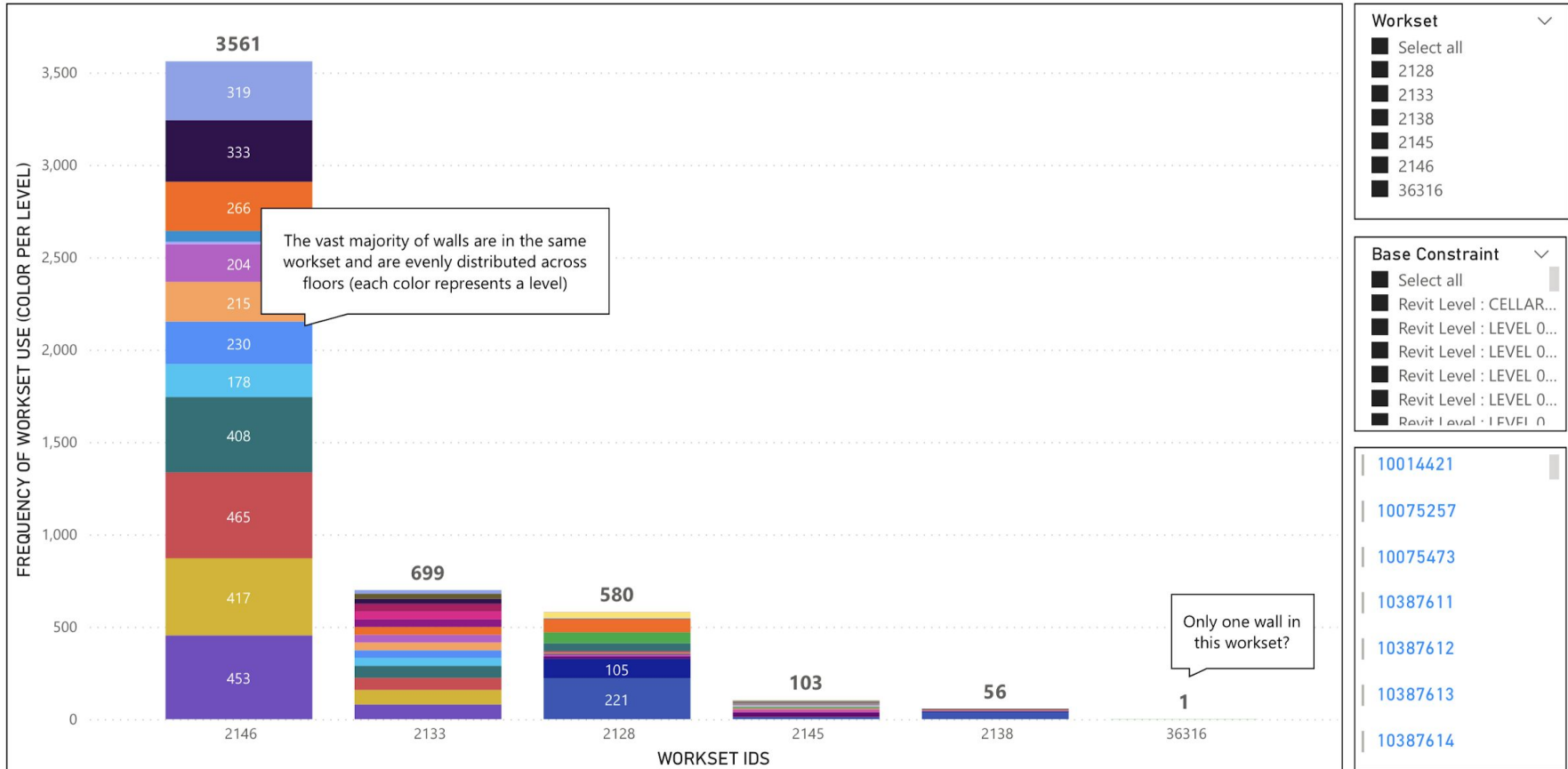| Name | Type | ID | Wall_Width | Category | Category | Curtain version | Locked | Range | Angle 2 | Angle 1 | Deletabl | Id | Family and Type | Family |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Concrete_Core 22" | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412 | 5398956 | 22 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5398956 | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412650 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 22" | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412 | 5398957 | 22 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5398957 | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412650 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 22" | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412 | 5398958 | 22 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5398958 | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412650 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 22" | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412 | 5398959 | 22 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5398959 | Revit Host Type : Basic Wall : Concrete_Core 22" : id 5412650 | Revit Host Type : Basic Wall : Concrete_C |
| HA Wall_S42.2 | Revit Host Type : Basic Wall : HA Wall_S42.2 : id 588669 | 5399158 | 5.25 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | FALSE | | | | TRUE | 5399158 | Revit Host Type : Basic Wall : HA Wall_S42.2 : id 588669 | Revit Host Type : Basic Wall : HA Wall_S42 |
| Concrete_Foundation 12" | Revit Host Type : Basic Wall : Concrete_Foundation 12" : i | 5413295 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | FALSE | | | | TRUE | 5413295 | Revit Host Type : Basic Wall : Concrete_Foundation 12" : id 308814 | Revit Host Type : Basic Wall : Concrete_F |
| Concrete_Foundation 12" | Revit Host Type : Basic Wall : Concrete_Foundation 12" : i | 5414040 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5414040 | Revit Host Type : Basic Wall : Concrete_Foundation 12" : id 308814 | Revit Host Type : Basic Wall : Concrete_F |
| Concrete_Foundation 12" | Revit Host Type : Basic Wall : Concrete_Foundation 12" : i | 5414487 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5414487 | Revit Host Type : Basic Wall : Concrete_Foundation 12" : id 308814 | Revit Host Type : Basic Wall : Concrete_F |
| Concrete_Foundation 12" | Revit Host Type : Basic Wall : Concrete_Foundation 12" : i | 5415156 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415156 | Revit Host Type : Basic Wall : Concrete_Foundation 12" : id 308814 | Revit Host Type : Basic Wall : Concrete_F |
| Concrete_Foundation 12" | Revit Host Type : Basic Wall : Concrete_Foundation 12" : i | 5415289 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415289 | Revit Host Type : Basic Wall : Concrete_Foundation 12" : id 308814 | Revit Host Type : Basic Wall : Concrete_F |
| Concrete_Foundation 12" | Revit Host Type : Basic Wall : Concrete_Foundation 12" : i | 5415418 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415418 | Revit Host Type : Basic Wall : Concrete_Foundation 12" : id 308814 | Revit Host Type : Basic Wall : Concrete_F |
| Concrete_Core 20" | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415 | 5415752 | 20 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415752 | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415705 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 20" | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415 | 5415811 | 20 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415811 | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415705 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 20" | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415 | 5415862 | 20 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415862 | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415705 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 20" | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415 | 5415962 | 20 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5415962 | Revit Host Type : Basic Wall : Concrete_Core 20" : id 5415705 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 18" | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416 | 5416103 | 18 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416103 | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416031 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 18" | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416 | 5416152 | 18 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416152 | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416031 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 18" | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416 | 5416204 | 18 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416204 | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416031 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 18" | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416 | 5416321 | 18 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416321 | Revit Host Type : Basic Wall : Concrete_Core 18" : id 5416031 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 12" | Revit Host Type : Basic Wall : Concrete_Core 12" : id 14398 | 5416433 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416433 | Revit Host Type : Basic Wall : Concrete_Core 12" : id 1493868 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 12" | Revit Host Type : Basic Wall : Concrete_Core 12" : id 14398 | 5416481 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416481 | Revit Host Type : Basic Wall : Concrete_Core 12" : id 1493868 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 12" | Revit Host Type : Basic Wall : Concrete_Core 12" : id 14398 | 5416538 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416538 | Revit Host Type : Basic Wall : Concrete_Core 12" : id 1493868 | Revit Host Type : Basic Wall : Concrete_C |
| Concrete_Core 12" | Revit Host Type : Basic Wall : Concrete_Core 12" : id 14398 | 5416596 | 12 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416596 | Revit Host Type : Basic Wall : Concrete_Core 12" : id 1493868 | Revit Host Type : Basic Wall : Concrete_C |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5416683 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5416683 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5417517 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5417517 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5417569 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5417569 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5418202 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5418202 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5418206 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5418206 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5418208 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5418208 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5423634 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5423634 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5423840 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5423840 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5423874 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5423874 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5424014 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5424014 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5426569 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5426569 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5426572 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5426572 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5426576 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5426576 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427052 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427052 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427094 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427094 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W34.2 | Revit Host Type : Basic Wall : HA Wall_W34.2 : id 5388616 | 5427189 | 6.125 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | FALSE | | | | TRUE | 5427189 | Revit Host Type : Basic Wall : HA Wall_W34.2 : id 588661 | Revit Host Type : Basic Wall : HA Wall_W3 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5427395 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427395 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5427431 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427431 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427689 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427689 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427747 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427747 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427820 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427820 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427838 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427838 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5427870 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5427870 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428121 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428121 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428254 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428254 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428293 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428293 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428411 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428411 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428619 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428619 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428708 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428708 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5428820 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5428820 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5429121 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429121 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5429279 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429279 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5429308 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429308 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5429350 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429350 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5429403 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429403 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5429481 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429481 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.2 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 ( | 5429510 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429510 | Revit Host Type : Basic Wall : HA Wall_W24.2 : id 588664 | Revit Host Type : Basic Wall : HA Wall_W2 |
| HA Wall_W24.1 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | 5429613 | 5 | Revit Category : Walls : id -2000011 | Revit Category : Walls : id -2000011 | | TRUE | | | | TRUE | 5429613 | Revit Host Type : Basic Wall : HA Wall_W24.1 : id 4780387 | Revit Host Type : Basic Wall : HA Wall_W2 |

Sheet1

# Outliers are bad for Machine Learning

Outliers: TOP OFFSET (each column represents a top offset value, the height of the column shows the qty of walls with that value)
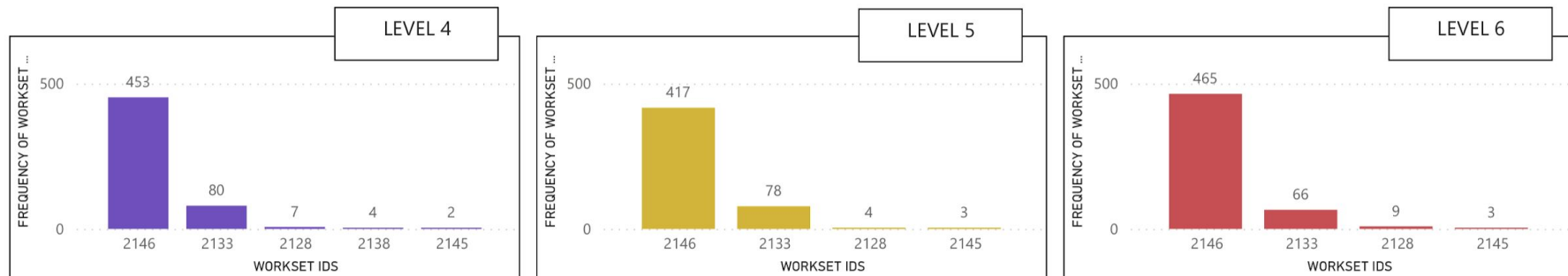
Outliers: WORKSET USAGE (each column represents a workset, height shows number of walls in that workset)
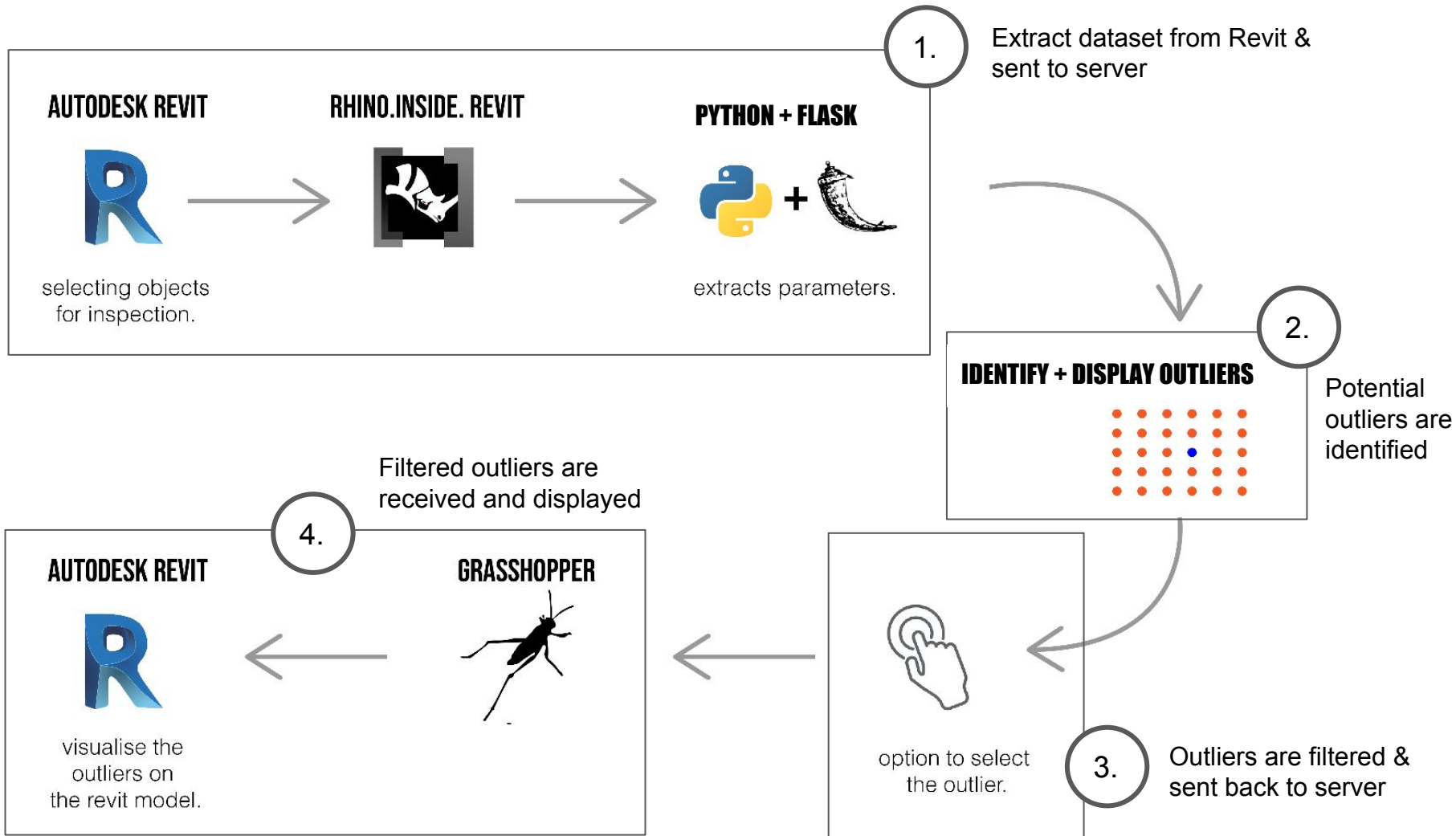
Outliers: NUMBER OF WORKSETS IN USE BY LEVEL (each column represents a workset)

the "four column type"

LEVEL 4

LEVEL 5

LEVEL 6

the "three column type"

LEVEL 3

LEVEL 2

LEVEL 8

**1.** Extract dataset from Revit & sent to server

AUTODESK REVIT

selecting objects for inspection.

RHINO.INSIDE. REVIT

PYTHON + FLASK

extracts parameters.

**2.** Potential outliers are identified

IDENTIFY + DISPLAY OUTLIERS

**3.** Outliers are filtered & sent back to server

option to select the outlier.

GRASSHOPPER

**4.** Filtered outliers are received and displayed

AUTODESK REVIT

visualise the outliers on the revit model.

# 1. Extract data from Revit

Revit Model

Rhino.Inside
(Use grasshopper with Revit)

Hops to Flask
(Send a stream of data to a
flask server)

# 2. Identify possible outliers

**Flask Server**
(Call server to access JSON)

↓

**SK Learn**
(Univariate outliers,
Mahalanobis, Box Plot,
Scatter Plot)

↓

**Flask Server**
(Send list of ids to server)

```python
#Detecting Mutivariate Outliers
#each parameter at the input should be numerical
def robust_mahalanobis_method(df):
    #Minimum covariance determinant
    rng = np.random.RandomState(0)
    real_cov = np.cov(df.values.T)
    X = rng.multivariate_normal(mean=np.mean(df, axis=0), cov=real_cov, size=506)
    cov = MinCovDet(random_state=0).fit(X)
    mcd = cov.covariance_ #robust covariance metric
    robust_mean = cov.location_  #robust mean
    inv_covmat = sp.linalg.inv(mcd) #inverse covariance metric

    #Robust M-Distance
    x_minus_mu = df - robust_mean
    left_term = np.dot(x_minus_mu, inv_covmat)
    mahal = np.dot(left_term, x_minus_mu.T)
    md = np.sqrt(mahal.diagonal())

    #Flag as outlier
    outlier = []
    C = np.sqrt(chi2.ppf((1-0.00001), df=df.shape[1]))#degrees of freedom = number of variables
    for index, value in enumerate(md):
        if value > C:
            outlier.append(index)
        else:
            continue
    outlier_id = df.index[outlier]
    return outlier_id
```

```python
#Outlier Functions...
#Detecting Univariate Outliers
def univariate_outlier(df,variable,q1_percent,q3_percent):
    q1 = df[variable].quantile(q1_percent)
    q3 = df[variable].quantile(q3_percent)
    iqr = q3-q1
    inner_fence = 1.5*iqr
    outer_fence = 3*iqr

    #inner fence lower and upper end
    inner_fence_le = q1-inner_fence
    inner_fence_ue = q3+inner_fence

    #outer fence lower and upper end
    outer_fence_le = q1-outer_fence
    outer_fence_ue = q3+outer_fence

    outliers_prob = []
    outliers_poss = []
    for i in range(len(df[variable])):
        x = df[variable].values[i]
        index = df[variable].take([i]).index[0]
        if x < outer_fence_le or x > outer_fence_ue:
            outliers_prob.append(index)

    for i in range(len(df[variable])):
        x = df[variable].values[i]
        index = df[variable].take([0]).index[0]
        if x <= inner_fence_le or x >= inner_fence_ue:
            outliers_poss.append(index)
    return outliers_prob, outliers_poss
```

# 3. Visualize, Explore, and Select Outliers

### Flask Server
(Call server to access JSON)

↓

## Web Dashboard
(React / Typescript /
Plotly / Material-UI)

↓

### Flask Server
(Send list of ids to server)

# 4. View outliers in Revit

**Flask Server**
(Call server to access JSON)

**Rhino.Inside / Hops**
(Get selected ids)

**Revit**
(Visualize selected elements)

# 4. View outliers in Revit

**Flask Server**
(Call server to access JSON)

↓

**Rhino.Inside / Hops**
(Get selected ids)

↓

**Revit**
(Visualize selected elements)

# Future Development

- **Use Hops instead of Excel for data export**
  Automate export/ analysis procedure

- **Automatic highlighting in Revit**
  Hops should return IDs into Grasshopper so the user does not need to copy & paste

- **"Remember" User's Input**
  Once the user has confirmed that a particular parameter on a particular element its not an outlier, that element/ parameter would not be flagged as an outlier again

- **Recipes**
  Learn from experience which parameters are likely to be worthwhile to cross-reference, and prioritize these for new projects.

- **Automatic recognition of meta-patterns**
  Implement multi-variable regressions

# Thanks!

https://github.com/jperaino/cuckoo

# PRESENTATION PITCH CONCEPT OUTLINE
## (for discussion)

Clearly Identify:

Key Idea & value proposition
    Real problem.  Own the "non-flashy"
    NOTHING existed before saturday
    Errors & Omissions are a real problem. (RFIs, litigation)
        Specific uses: data consistency in naming, wrong family use,     i.
.        Incorrect size
    We understand not all "outliers" are mistakes.  Filtering is key (enter ML)
Key tech innovations used (NEED TO HIT THIS… what are they????)
    Identify open source libraries used
    Any ML  usage? (PYOD)
Key future potential
    Model training from multiple models

Key Idea & value proposition
Key future potential

**MB**

**1.   MAIN VISION MAIN INTRO**
Introduce team, main idea, vision, base concepts, etc

**2. VALUE POTENTIAL**
Base proof of concept: pbi diagrams

**AO**

**3. HOW IT IS DONE**
Data flow diagram
-this is key. Show whole picture, understanding

**SSP**

Key Idea & value proposition
Key tech innovations used
Key future potential

**6. FINAL SELL PITCH**
Close
RECAP, VISION, FUTURE POTENTIAL

**MB**

**5. TECHNOLOGICAL INNOVATION Our Hack…**
processes & results

**4. TECHNOLOGICAL INNOVATION**
Our Hack…
processes & results

**JP**

**ZH**

Key Idea & value proposition
Key tech innovations used
Key future potential

Key tech innovations used