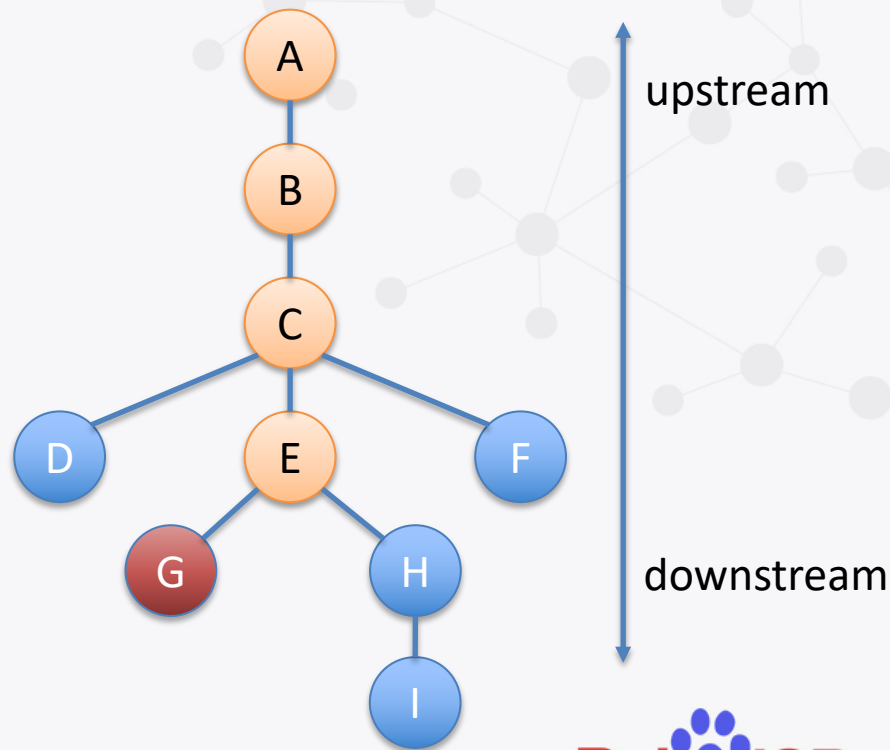# Service Diagnosis

- Quick diagnosis upon service faults
- Narrow down the checking scope
  - Perform stop-loss actions
    - Traffic switching, release rollback
  - Understand how the fault emerges
- Automatic metric screening
  - Reduce diagnosis time by 80%
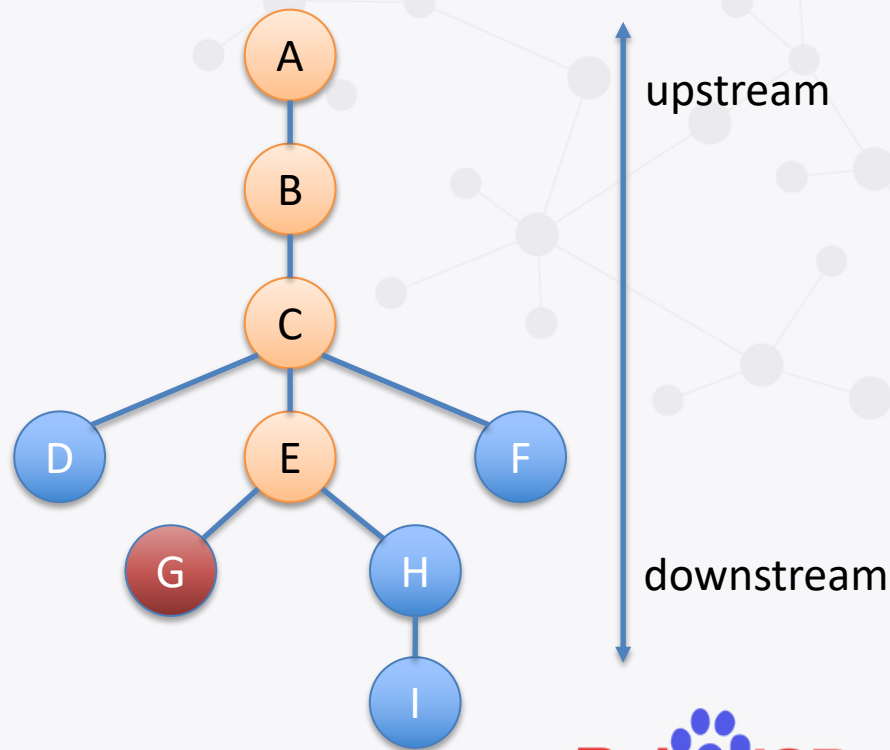  - Minimum manual configuration

# Diagnosing a Fault

- PV lost perceived at front end
- Get the call graph
  - Modules as nodes
  - Calls as edges
- Examine the modules
  - upstream to downstream
  - Look at certain metrics
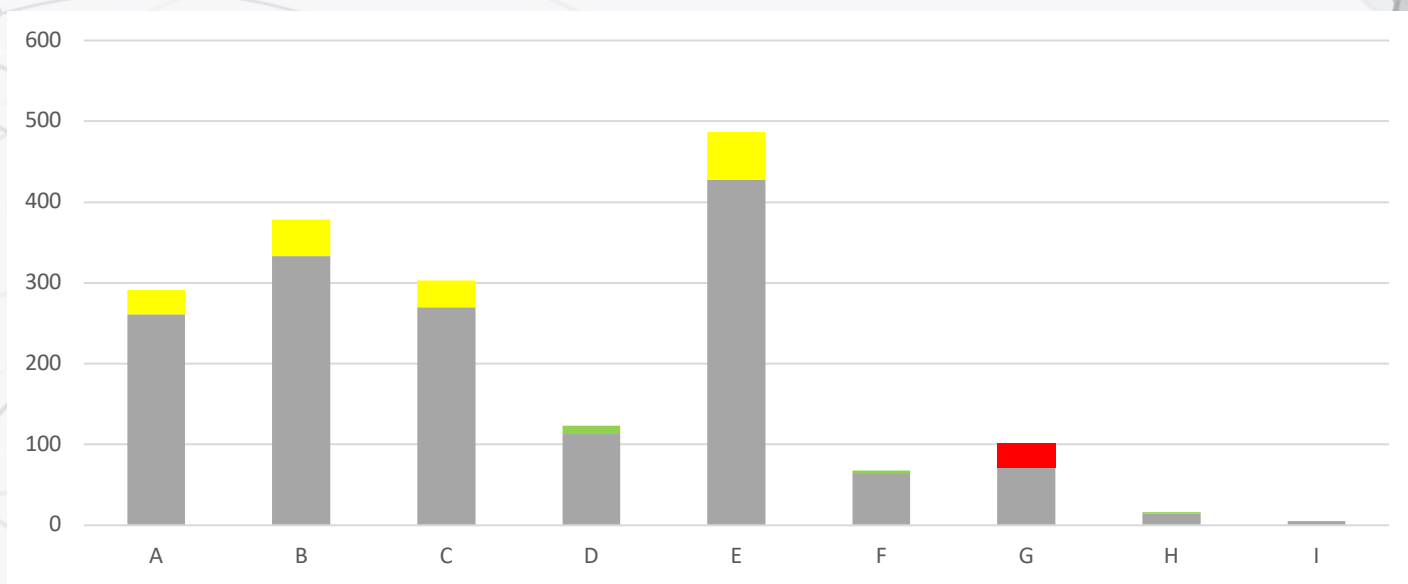  - Decide the root cause module

upstream

downstream

Bai du IOP

# Diagnosis Tool

- Show the call graph
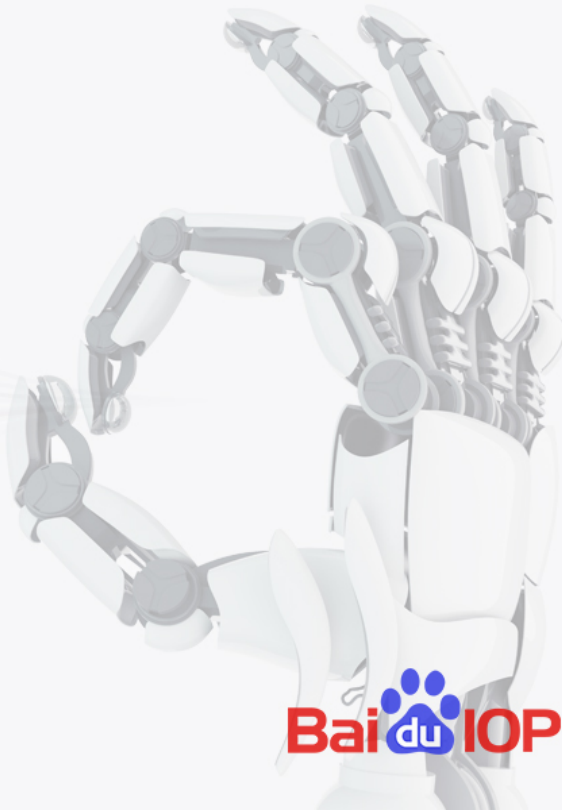- Mark the abnormal modules
  - How to mark?

# Module Abnormal Detection

- Golden metric, e.g. response time

# Semi-automatic Approach

- Synchronize call graphs
  - Continuous call graph evolution
  - Standard RPC middleware
- Golden metric configuration
  - New fault type requires new metrics
- Multiple call graphs
  - One call graph per datacenter

Bai du IOP

# Automatic Metric Screening

- No dependency on golden metrics
  - Check all modules, all metrics
- Find out abnormal metrics
  - Flexible to metric diversity
- Recommendation
  - Quantitative anomaly measurement
  - Legible

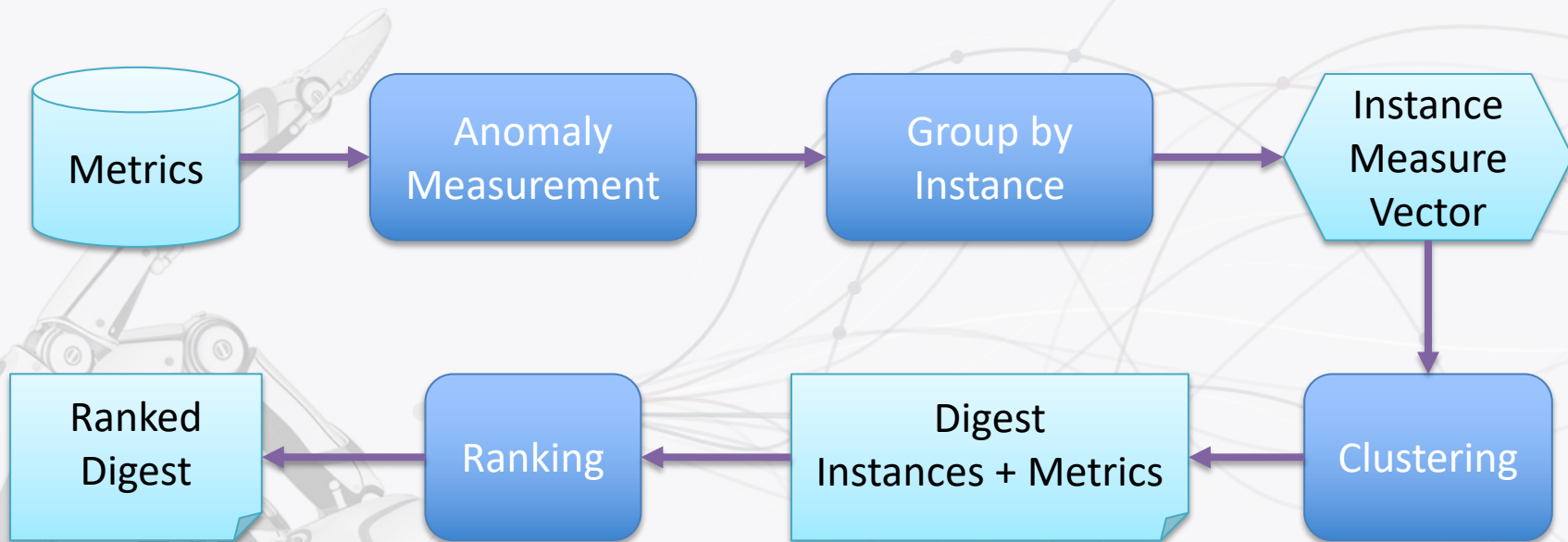# Screening Instance Level Metrics

Metrics → Anomaly Measurement → Group by Instance → Instance Measure Vector

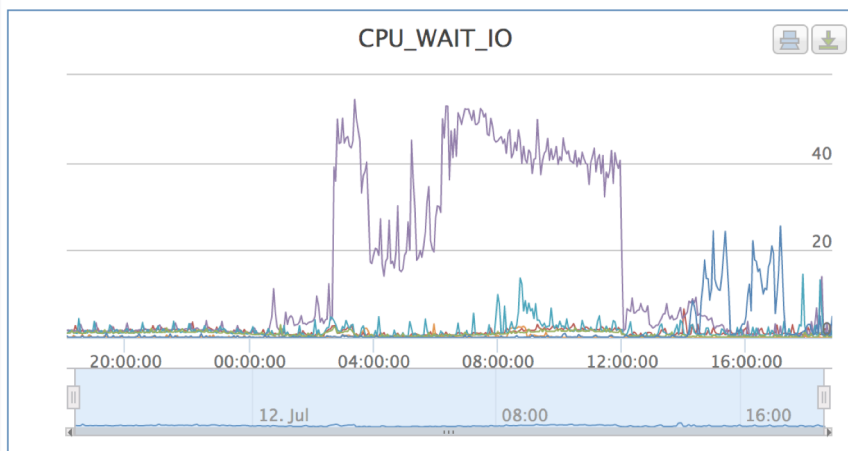Instance Measure Vector → Clustering → Digest Instances + Metrics → Ranking → Ranked Digest
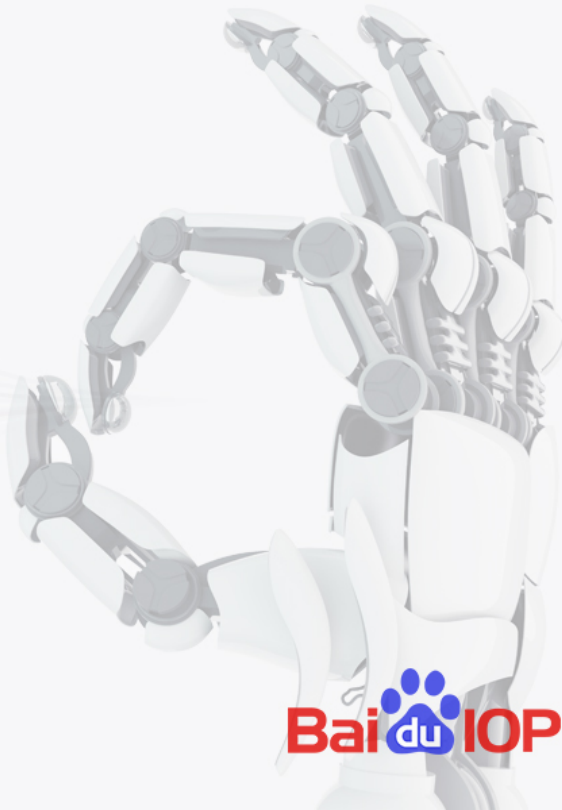
# Anomaly Measurement (I)

- Sudden changes
  - Many faults reflect on sudden changes on certain metrics
  - Easy to understand
  - Easy to detect
    - Compare data after faults against those before faults



Bai du IOP

# Anomaly Measurement (II)

- Data
  - 60 mins before the fault $\{x_i\}$
  - 5 mins after the fault $\{y_j\}$

- Measure
  - Conditional probability $P(\{y_j\}|\{x_i\})$
  - Single point probability
    - Overflow $P(y \geq y_j|\{x_i\})$
    - Underflow $P(y \leq y_j|\{x_i\})$

# Anomaly Measurement (III)

- Kernel density estimation
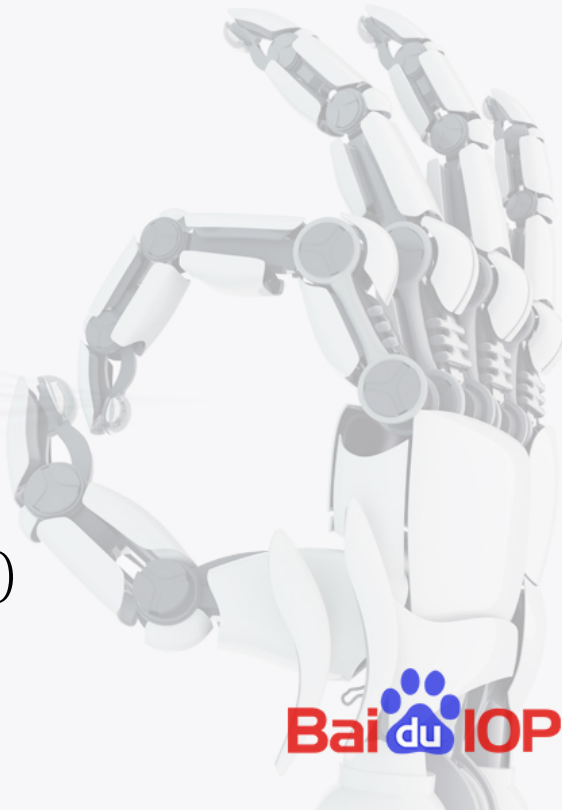  - $f(y) = \frac{1}{n}\sum_{i=1}^{n} \mathcal{N}(y; x_i, \sigma)$

  - $P\left(y \geq y_j \middle| \{x_i\}\right) = \int_{y_j}^{+\infty} f(y)$

  - $P\left(y \leq y_j \middle| \{x_i\}\right) = \int_{-\infty}^{y_j} f(y)$
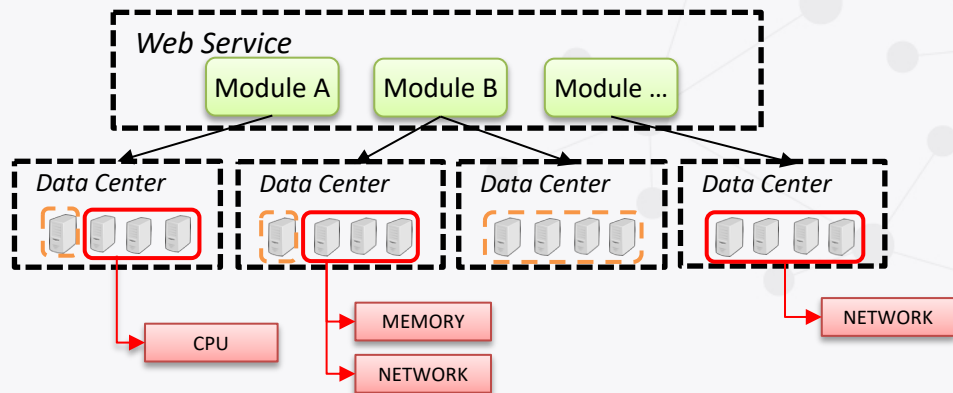
  - Kernels: Gaussian, Poisson, Beta
- Use log probability to combine
  - $\log P_{\text{overflow}}\left(\{y_j\} \middle| \{x_i\}\right) = \frac{1}{m}\sum_{j=1}^{m} \log P_{\text{overflow}}\left(y \geq y_j \middle| \{x_i\}\right)$

Bai du IOP

# Clustering

- Anomaly measure vector
  - $D_m = -\log P(\{y_j\}|\{x_i\})$ for metric $m$
  - $\langle D_1, D_2, \ldots \rangle$ for each instance
- Clustering
  - Instances within a module×DC
  - DBSCAN
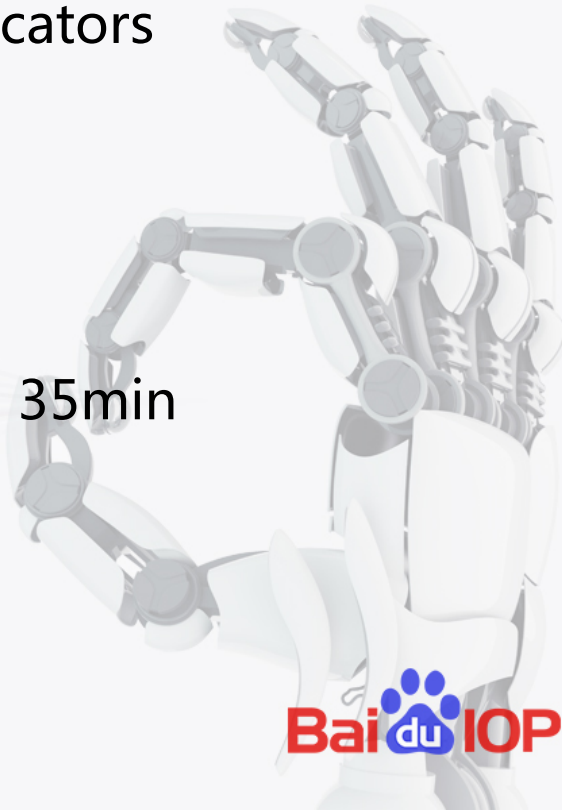  - Digest = instance set + abnormal metric set

# Ranking the Digests

- Proportion of the instances
- Number of abnormal metrics
- Anomaly degree
  - $D_m = -\log P(\{y_j\}|\{x_i\})$
- Train a ranker

| Module | Portion | Instances | Metrics | |
|--------|---------|-----------|---------|---|
| G.DC1 | 0.3 | 1.G.DC1<br>2.G.DC1<br>3.G.DC1<br>... | CPU_SERVER_LOADAVG_1<br>NET_TCP_ACTIVE_OPENS<br>CPU_IDLE<br>CPU_HT_IDLE<br>NET_TCP_TIME_WAIT<br>CPU_INTERRUPT<br>NET_TCP_OUT_SEGS<br>CPU_SERVER_LOADAVG_5<br>NET_TCP_IN_SEGS<br>... | 22.801403<br>16.90959<br>15.137489<br>14.600515<br>11.838425<br>11.682794<br>11.613833<br>11.550683<br>11.136453 |
| E.DC1 | 0.1 | 7.E.DC1<br>9.E.DC1<br>... | NET_TCP_CURR_ESTAB<br>NET_TOTAL_SOCKETS_USED<br>DISK_TOTAL_READ_REQ<br>... | 15.696724<br>14.688062<br>11.241439 |

Bai du IOP

# Offline Evaluation Result

- Use Linux standard system performance indicators
- 70 historic cases
  - Root cause can reflect on Linux indicators
- 60 ranked root cause digest as top 1

- Human diagnosis time: 6min~152min, mean 35min
- Algorithm execution time: <=6min

Bai du IOP

# THANK YOU !

BAIDU@IOP出品