# *A Case for Packet Sampling*

**FOKUS**

**Fraunhofer** Institute for Open
Communication Systems

Tanja Zseby, zseby@fokus.fhg.de

Competence Center for Autonomic Networking Technologies

# *Motivation: FloCon 2005*

FloCon05 participants:


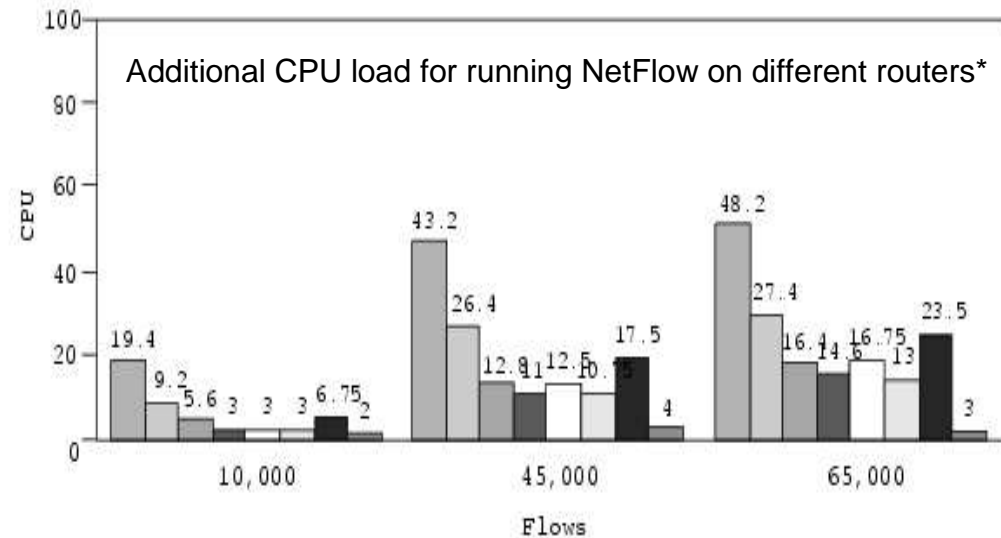"We don't believe in Sampling"


- Happy to use flow data

- Very skeptical to packet sampling

FOKUS

Fraunhofer Institute for Open
Communication Systems

# *The Problem: Limited Resources*

- Full packet capture at each node not feasible
  - Increasing data rates
  - Hardware costs
  - Privacy concerns
- Resources are limited
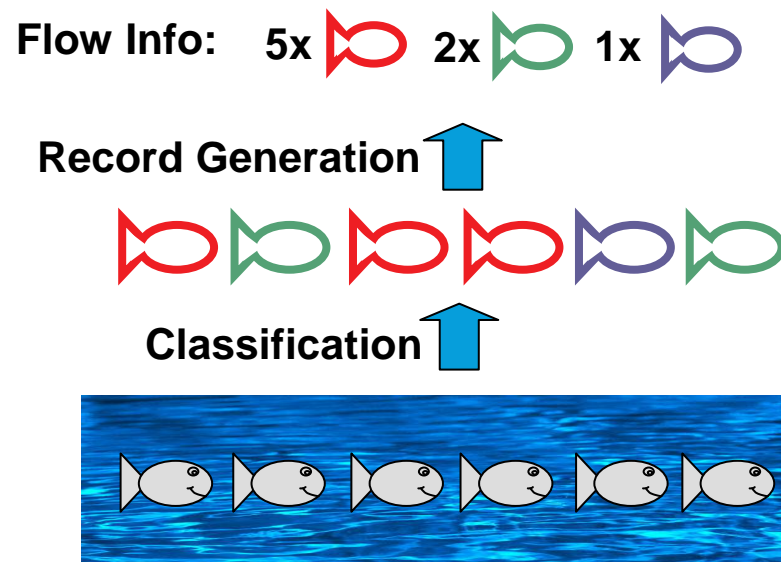  - Storage
  - Processing
  - Transmission



Additional CPU load for running NetFlow on different routers*

**We cannot measure everything**

*source: NetFlow Performance Analysis, Cisco white paper
http://www.cisco.com/warp/public/cc/pd/iosw/prodlit/ntfo_wpa.jpg

**FOKUS**

Fraunhofer Institute for Open Communication Systems

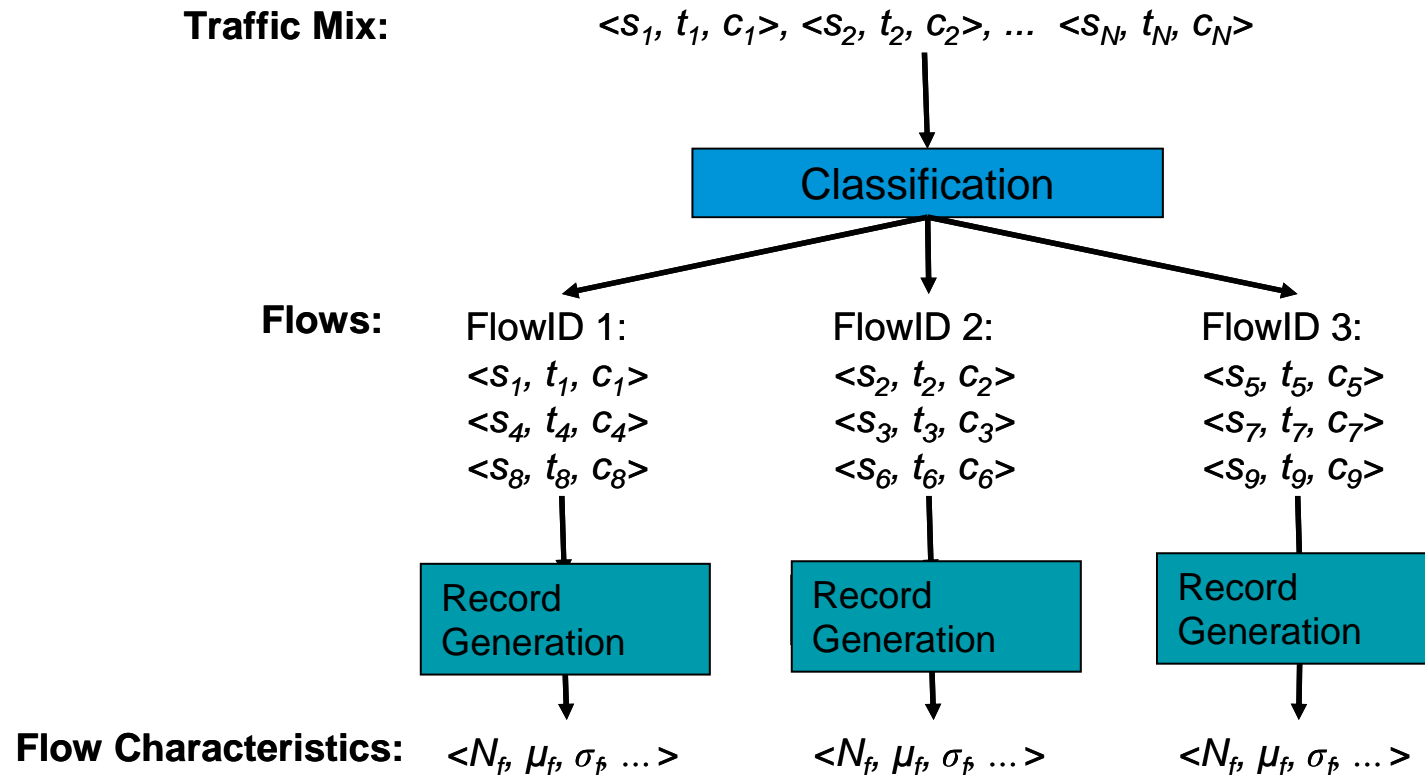# Solution1: Flow Data

- Grouping of packets into flows (classification)

- Reporting of flow information only

- Disadvantages:
  - Per-packet information is lost
  - Information and effort depends on flow definition

Flow Info:  5x  2x  1x

Record Generation

Classification

# *Flow Data Generation*

**Traffic Mix:** $\langle s_1, t_1, c_1 \rangle, \langle s_2, t_2, c_2 \rangle, \ldots \langle s_N, t_N, c_N \rangle$

Classification

**Flows:**

FlowID 1:
$\langle s_1, t_1, c_1 \rangle$
$\langle s_4, t_4, c_4 \rangle$
$\langle s_8, t_8, c_8 \rangle$

FlowID 2:
$\langle s_2, t_2, c_2 \rangle$
$\langle s_3, t_3, c_3 \rangle$
$\langle s_6, t_6, c_6 \rangle$

FlowID 3:
$\langle s_5, t_5, c_5 \rangle$
$\langle s_7, t_7, c_7 \rangle$
$\langle s_9, t_9, c_9 \rangle$

Record Generation | Record Generation | Record Generation

**Flow Characteristics:** $\langle N_f, \mu_f, \sigma_f \ldots \rangle$  $\langle N_f, \mu_f, \sigma_f \ldots \rangle$  $\langle N_f, \mu_f, \sigma_f \ldots \rangle$
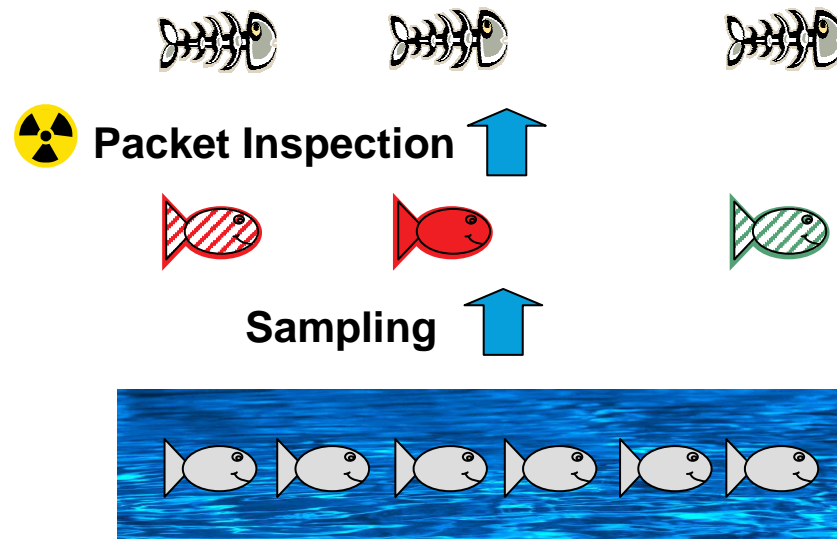
- Information about packets is discarded
- Available information depends on
  - Flow definition
  - Flow characteristics that are reported

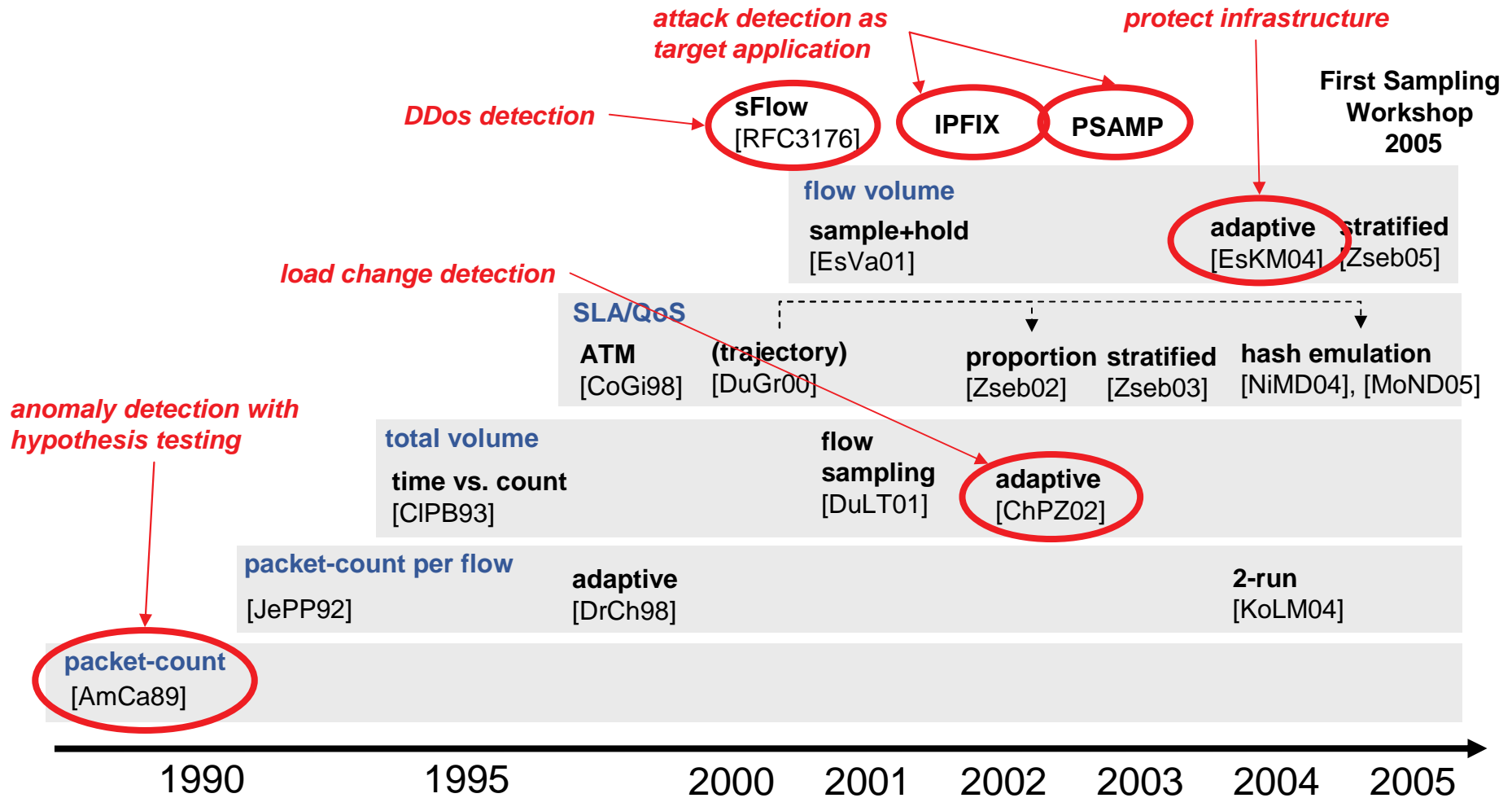FOKUS

Fraunhofer Institute for Open Communication Systems

# Solution2: Packet Sampling

- Random Selection of some packets
  - Report parts or full packet information
  - Estimation of metrics based on sample

- Provides different viewpoint
  - Packet data can reveal further information
  - Sampled data sufficient for some metrics

- Helps to protect measurement infrastructure during attack

**Packet Inspection**

**Sampling**

# *Sampling: State of Art*



**attack detection as target application**

**protect infrastructure**

**DDos detection**

sFlow [RFC3176]

IPFIX

PSAMP

**First Sampling Workshop 2005**

**flow volume**
sample+hold [EsVa01]

adaptive [EsKM04]

stratified [Zseb05]

**load change detection**

**SLA/QoS**
ATM [CoGi98]   (trajectory) [DuGr00]

proportion [Zseb02]   stratified [Zseb03]   hash emulation [NiMD04], [MoND05]

**anomaly detection with hypothesis testing**

**total volume**
time vs. count [ClPB93]

flow sampling [DuLT01]

adaptive [ChPZ02]

**packet-count per flow**
[JePP92]

adaptive [DrCh98]

2-run [KoLM04]

**packet-count**
[AmCa89]

1990   1995   2000   2001   2002   2003   2004   2005

# *Packet Sampling*

Real metric substituted by estimate

➔ Accuracy statement is essential

Accuracy depends on

- – Sampling scheme

- – Estimation method

- – Position of sampling process in measurement sequence

- – Population characteristics (e.g. variance of metric of interest)

Fraunhofer Institute for Open Communication Systems

# A Simple Example

**Goal: Estimation of packet proportions (e.g. TCP-SYN packets in a flow)**

Real proportion: $P = \dfrac{M}{N}$ Estimate: $\hat{P} = \dfrac{m}{n}$

Estimation Accuracy (random n-of-N): $\sigma_{\hat{P}} = \sqrt{\dfrac{P \cdot (1-P)}{n}} \cdot \sqrt{\dfrac{N-n}{N-1}}$

Confidence Limits: $Prob\left(\hat{P} - z_c \cdot \sigma_{\hat{P}} \le P \le \hat{P} + z_c \cdot \sigma_{\hat{P}}\right) = 1 - \alpha$

Example: - Measurement interval with N=10,000 packets
  - Random packet selection 1% (n=100)

$\hat{P} = 0.9$ ➜ $\sigma_{\hat{P}} = 0.03$ ➜ 0.8226 ≤ P ≤ 0.977, with 99% confidence

$\hat{P} = 0.1$ ➜ same accuracy

$\hat{P} = 0.5$ (worst case) ➜ $\sigma_{\hat{P}} = 0.05$ ➜ 0.371 ≤ P ≤ 0.629, with 99% confidence

**Works with other packet properties, too!**

FOKUS

Fraunhofer Institute for Open Communication Systems

# *Advise*

- Don't restrict your analysis to flow data
  - Include further viewpoints
  - Use sampling in addition or as alternative to flow data

- Trust the power of statistics
  - It's a mature and well established field
  - ➔ full range of proven techniques

- Use sampling where applicable
  - Applicability depends on traffic profile, metric of interest, accuracy demand
    - Sampled data sufficient to detect large events (high volumes, high packet counts)
    - May be sufficient to estimate #pkts with specific properties (e.g. SYN, VoIP packets, small packets, packets with same content, etc.)
    - Others ➔ depends on scenario
  - Difficulties with rare events (stealth attacks, slow port scans)
  - Not suitable to re-assemble connections (but filtering may be)

FOKUS

Fraunhofer Institute for Open Communication Systems

# Thank you for your attention!

Fraunhofer Institute for Open Communication Systems

FOKUS