# Intro to Stream Processing

Motivation
○○○○

Data Streams
○○○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○○

1

# Data Processing so far ...



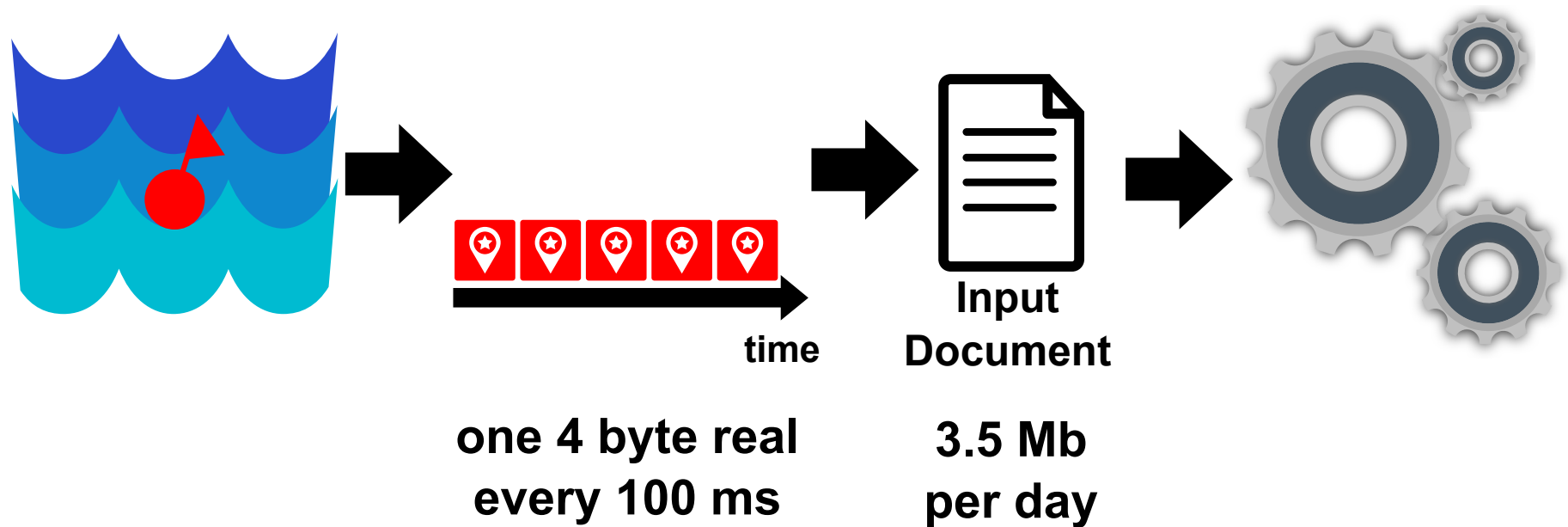**Input Document**  **Output Document**

Motivation
●○○○

Data Streams
○○○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○

3

# Sensor Data Example



**one 4 byte real per hour**

**Input Document**

**96 bytes per day**

Motivation
○●○○

Data Streams
○○○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○○

4

# Sensor Data Example



**Input Document**

one 4 byte real
every 100 ms

3.5 Mb
per day

**Motivation**
○●○○

Data Streams
○○○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○○○

5

# Sensor Data Example



**one million 4 byte reals every 100 ms**

**Input Document**

**3.5 Tb per day**

# Sensor Data Example
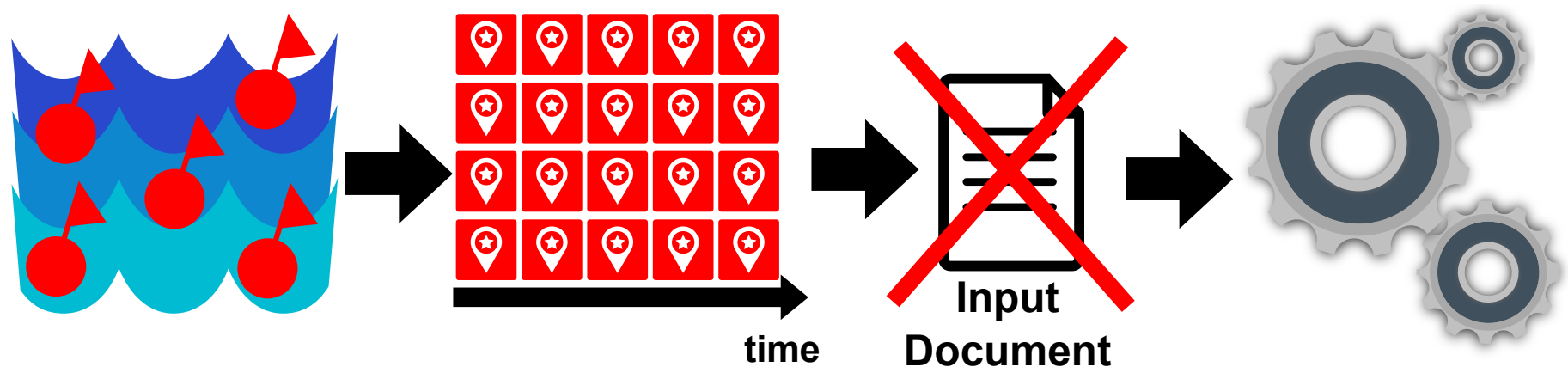
Stream of large unbounded data

too large for memory

too high latency for disk

We need real time processing!

Motivation
○○●○

Data Streams
○○○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○

7

# Sensor Data Example



**time**

**Input Document**

Process data stream directly

# Data Streams

Motivation
○○○○

**Data Streams**
●○○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○

9

# What is a Data Stream?

**Definition (Golab and Ozsu, 2003**

A data stream is a real-time, continuous, ordered (implicitly by arrival time of explicitly by timestamp) sequence of items. It is impossible to control the order in which items arrive, nor it is feasible to locally store a stream in its entirety.

# What is a Data Stream?

> **Definition (Golab and Ozsu, 2003**
>
> A data stream is a real-time, continuous, ordered (implicitly by arrival time of explicitly by timestamp) sequence of items. It is impossible to control the order in which items arrive, nor it is feasible to locally store a stream in its entirety.

- continous and sequential input
- typically unpredictable input rate
- can be large amounts of data
- not error free

Motivation
○○○○

Data Streams
○●○○○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○

11

# Data Stream Applications

- Online, real time processing
- Event detection and reaction
- Aggregation
- Approximation

Motivation
OOOO

**Data Streams**
OO●OOOOO

Reservoir Sampling
OOOOOOOOOOOOOOO

12

# Data Stream Example

Stock monitoring

Motivation
○○○○

Data Streams
○○○●○○○○

Reservoir Sampling
○○○○○○○○○○○○○○○

13

# Data Stream Example

Stock monitoring

Website traffic monitoring

Motivation
○○○○

Data Streams
○○○●○○○○

Reservoir Sampling
○○○○○○○○○○○○○○

14

# Data Stream Example

Stock monitoring

Website traffic monitoring

Network management

Motivation
○○○○

**Data Streams**
○○○●○○○○

Reservoir Sampling
○○○○○○○○○○○○○○

15

# Data Stream Example

Stock monitoring

Website traffic monitoring

Network management

Highway traffic

Motivation
○○○○

Data Streams
○○○●○○○○

Reservoir Sampling
○○○○○○○○○○○○○○

16

# Data Stream Characteristics



items

stream

Motivation
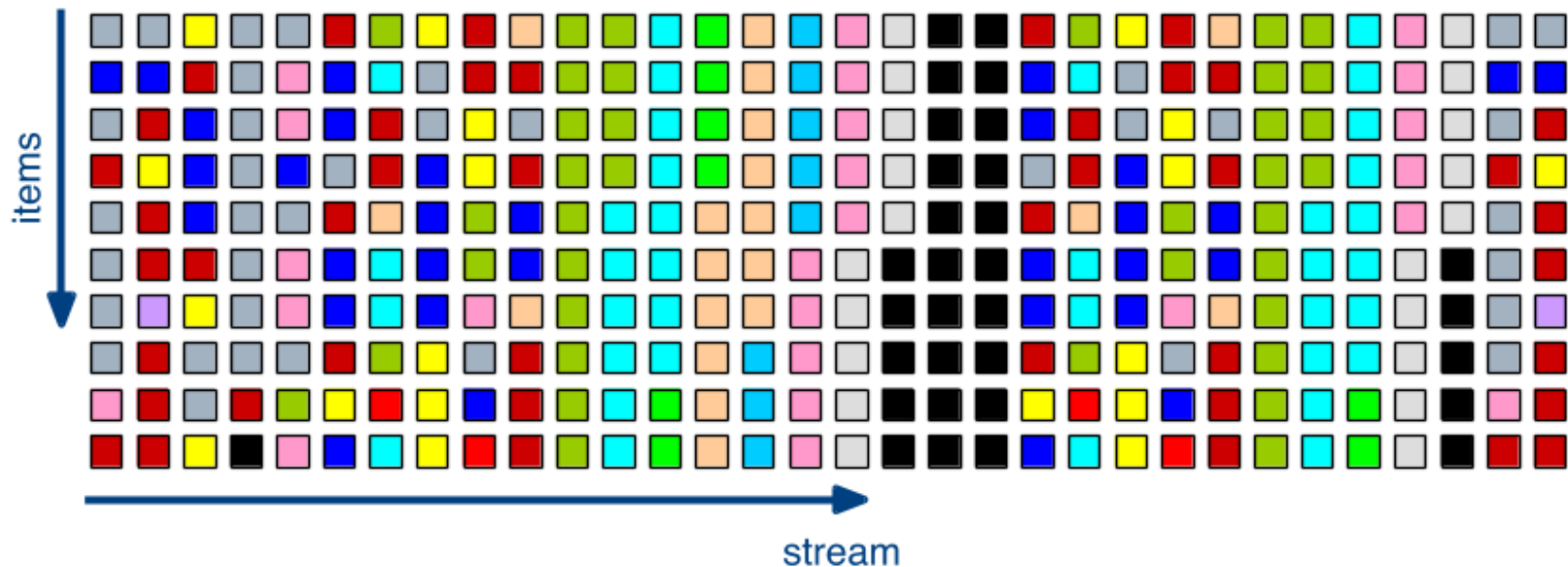○○○○

Data Streams
○○○○○●○○○

Reservoir Sampling
○○○○○○○○○○○○○○○

17

# Data Stream Characteristics



- All items have the same structure. For example a tuple or object: (sender, recipient, text body)

Motivation
○○○○

**Data Streams**
○○○○●○○○

Reservoir Sampling
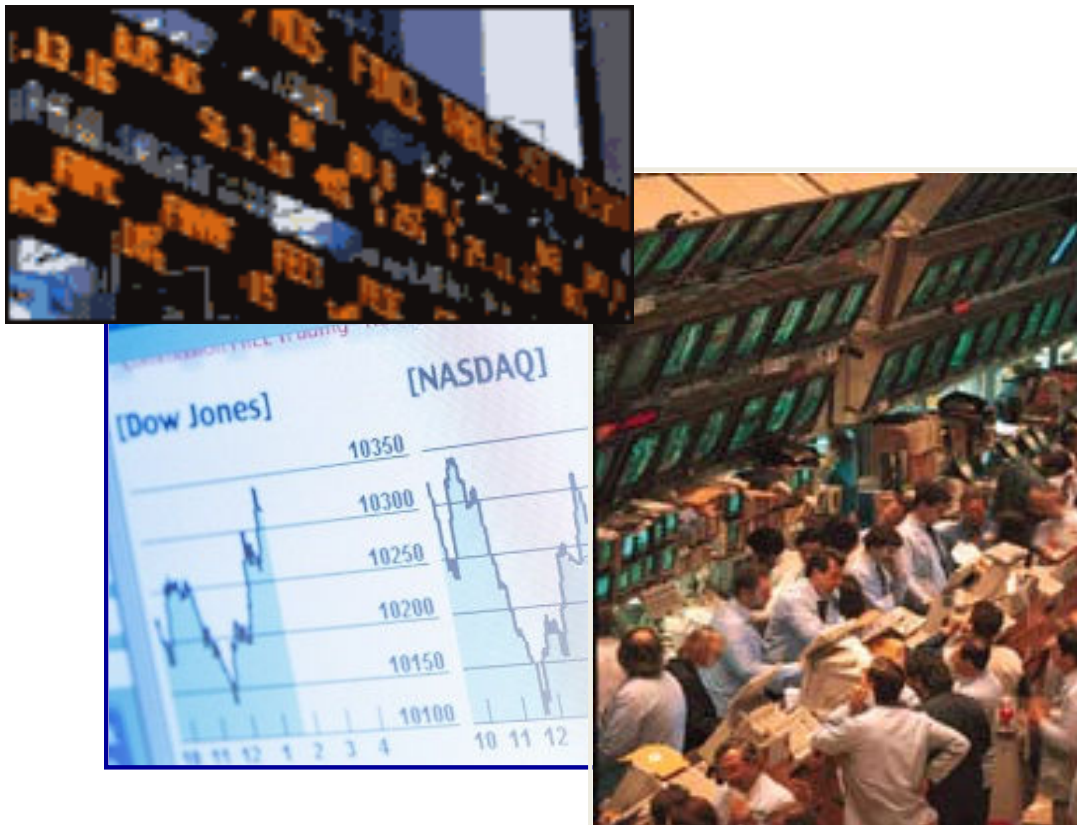○○○○○○○○○○○○○○○

18

# Data Stream Characteristics



- All items have the same structure. For example a tuple or object: (sender, recipient, text body)
- timestamps: explicite vs. implicite, physical vs. logical

Motivation
○○○○

**Data Streams**
○○○○●○○○

Reservoir Sampling
○○○○○○○○○○○○○○○○

19

# Data Streams

- Continuous sequences of data elements that are typically:
  - **Push-based** (like in publish/subscribe systems)
  - **Ordered** (e.g., by arrival time, or by explicit timestamps)
  - **Rapid** (e.g., ~ millions of messages/sec in market data)
  - **Potentially unbounded** (may have no (known) end)
  - **Time-sensitive** (real-time events, latency-critical)
  - **Time-varying** (in content and speed)
  - **Unpredictable** (autonomous data sources)

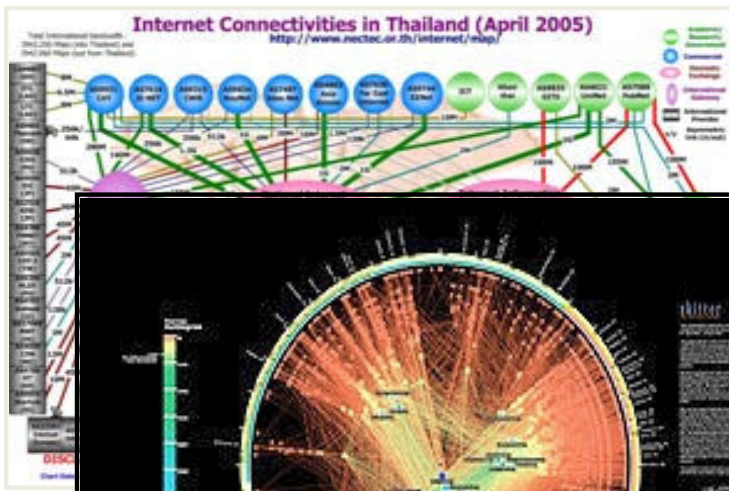# Example Applications

- Financial Services

**Example:**
- Trades(time, symbol, price, volume)

**Typical Applications:**
- Algorithmic Trading
- Foreign Exchange
- Fraud Detection
- Compliance Checking
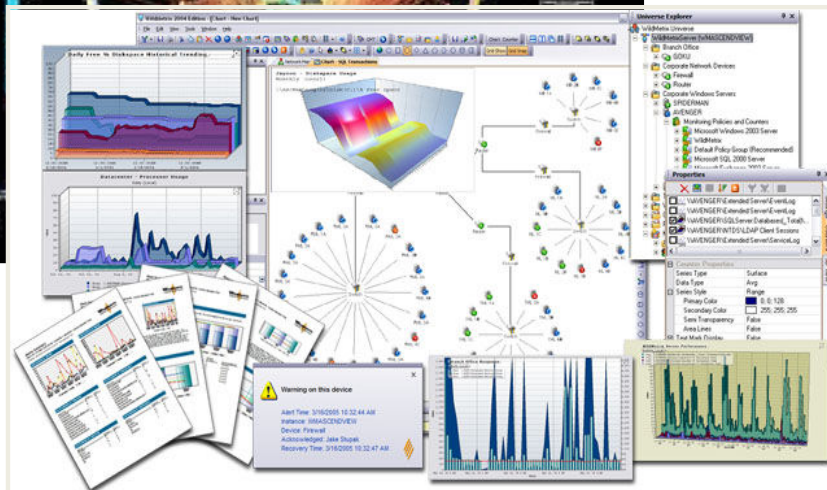
# Example Applications

- System and Network Monitoring



**Example:**
- Connections(time, srcIP, destIP, destPort, status)

**Typical Applications:**
- Server load monitoring
- Network traffic monitoring
- Detecting security attacks
  - Denial of Service
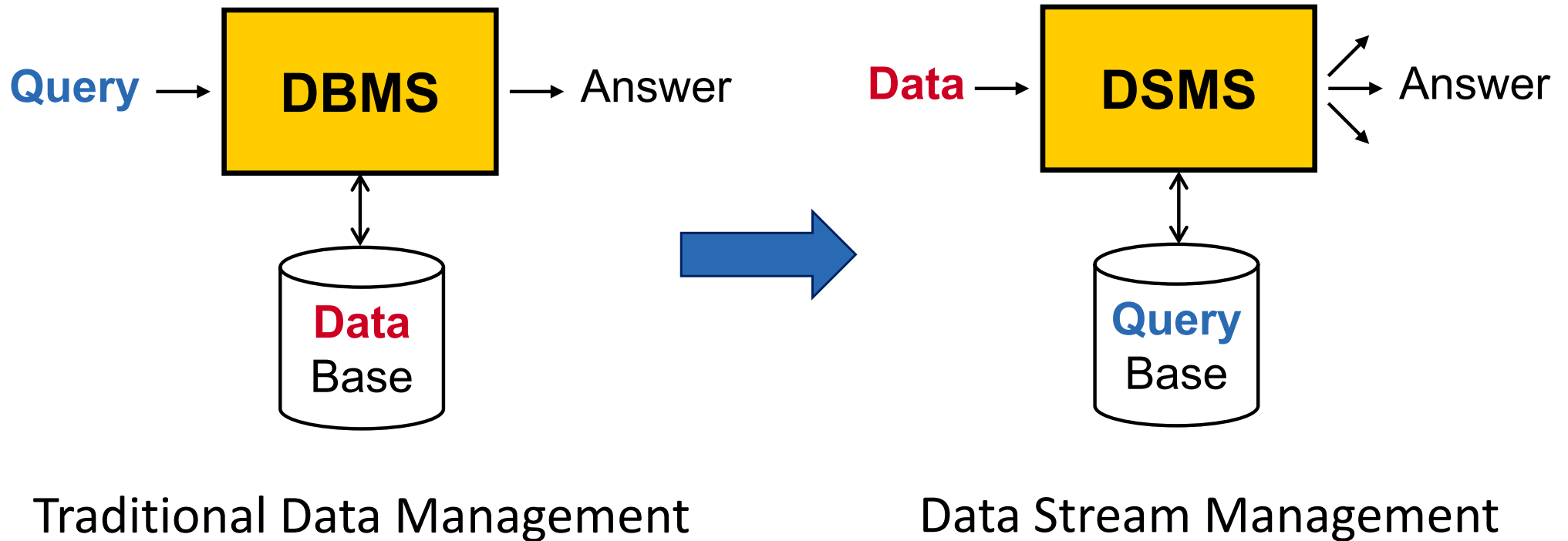  - Intrusion

# Example Applications

- ## Sensor-based Monitoring



**Example:**
- CarPositions(time, id, speed, position)

**Typical Applications:**
- Monitoring congested roads
- Route planning
- Rule violations
- Tolling

# A Paradigm Shift in Data Processing Model



Traditional Data Management

Data Stream Management

# DBMS   vs.   DSMS

- Persistent relations

- Read-intensive


- One-time queries



- Random access

- Access plan determined by query processor and physical DB design

- Transient streams

- Update-intensive (mostly append-only)

- Continuous queries (a.k.a., long-running, standing, or persistent queries)

- Sequential access

- Unpredictable data characteristics and arrival patterns