




Webinar Series 1 - 9
Magister Terapan PENS 2020

 Stream On
[Youtube.com/penstv](https://www.youtube.com/penstv)

Practical Data Mining

Case study : Intrusion Detection System

Iwan Syarif, MKom. MSc. PhD.
25 Agustus 2020

Search

This site University

ECS Home 

Who we are

Achievements

ECS people

Facilities

History of ECS

- Life & work of Eric Zepler
- 1957 - 1963
- 1963 - 1973
- 1963 - 1974

Our approach

Our University

Promoting Women in
Science and Engineering

University Home 

Electronics and Computer Science (ECS)

Home >

Iwan Syarif

Overview

Publications

ECS, Faculty of Physical Sciences and Engineering
University of Southampton
Southampton, United Kingdom. SO17 1BJ

Position: Postgraduate, submitted in Web and Internet
Science

Email: is1e08@ecs.soton.ac.uk

URI: <http://id.ecs.soton.ac.uk/person/24177> [browse]

Interests: artificial intelligence, computer network, data
mining, intrusion detection system, machine learning,
network security



Iwan Syarif

Qualifications

Master of Computer Science, University of Southampton,
UK, 2009

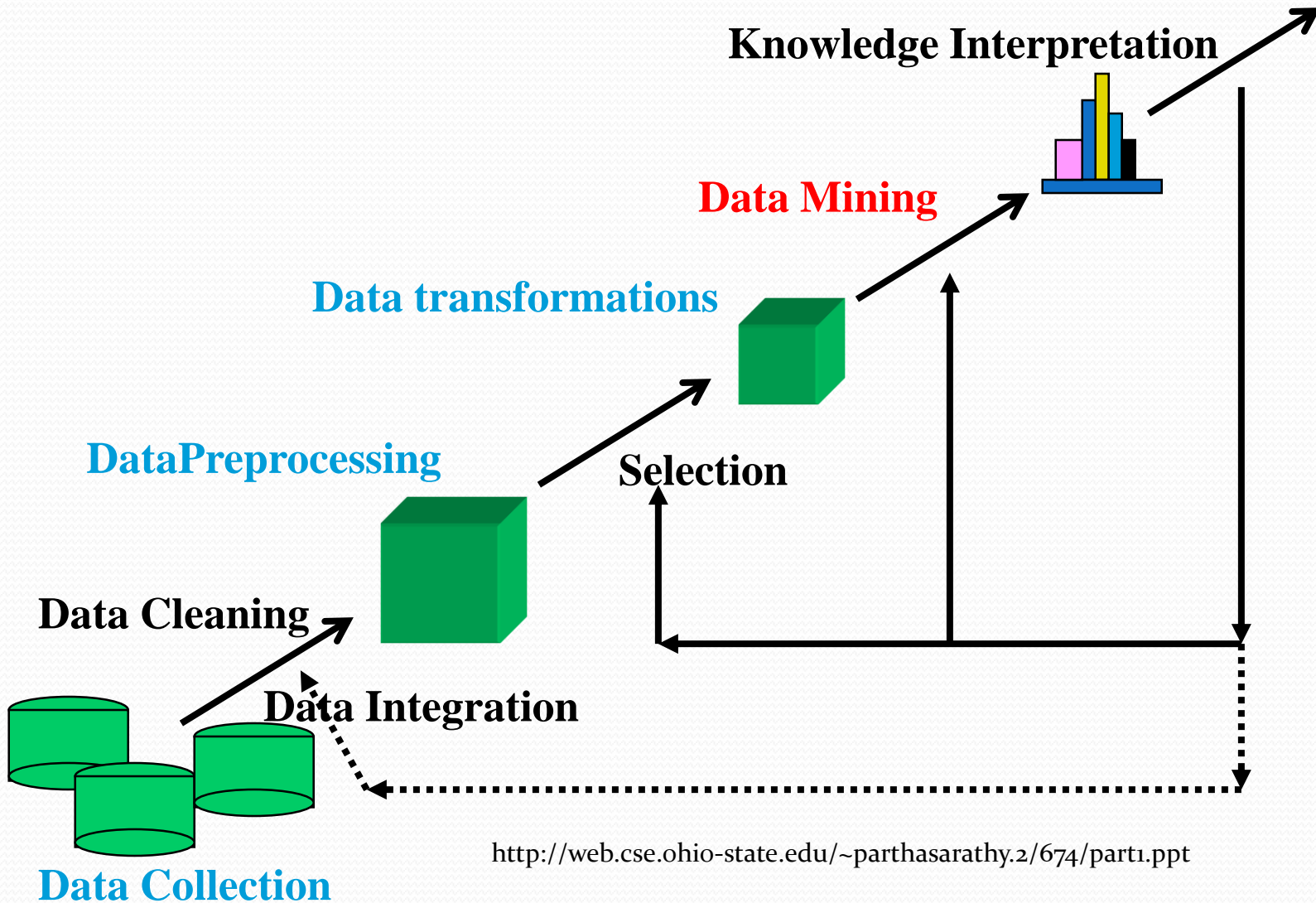
Master of Informatics, ITS Surabaya, Indonesia, 2003

Bachelor of Computer Engineering, ITS Surabaya, Indonesia, 1994

Data Mining

Data mining:
the core of knowledge discovery process.

Knowledge



Data Mining Approaches for Network Intrusion Detection System

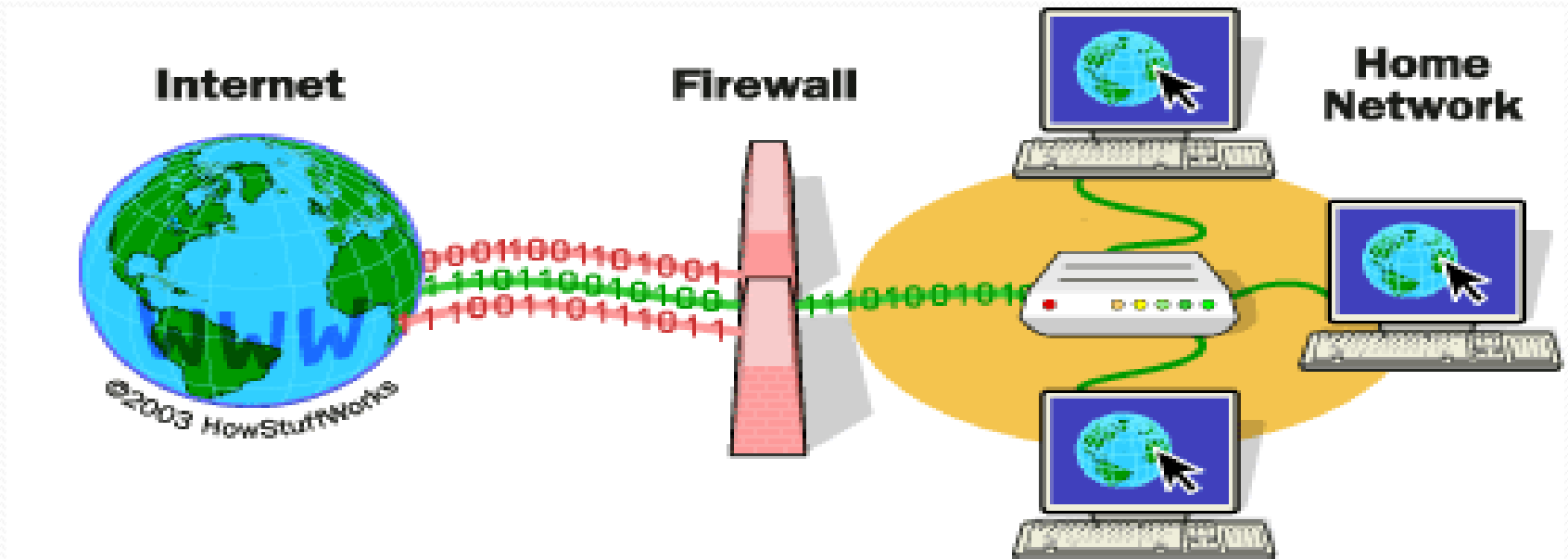


<https://norse-corp.com/map/>

Firewalls

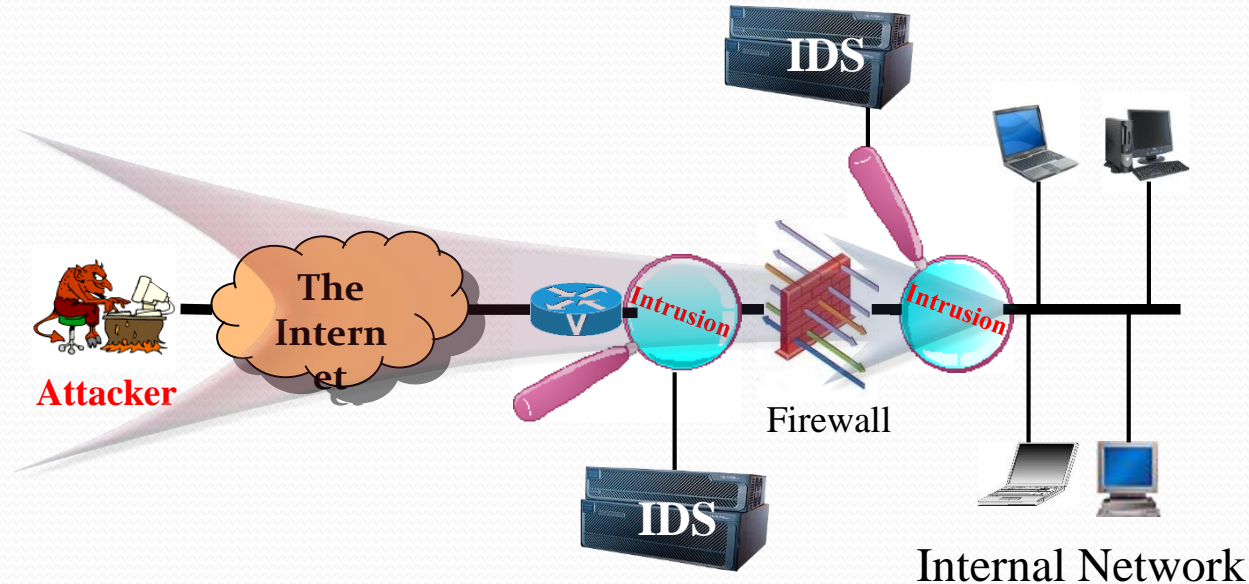


- Used to filter packets based on a combination of features
 - These are called packet filtering firewalls
 - Ex. Drop packets with destination port of 23 (Telnet)
- But why don't we just turn Telnet off?



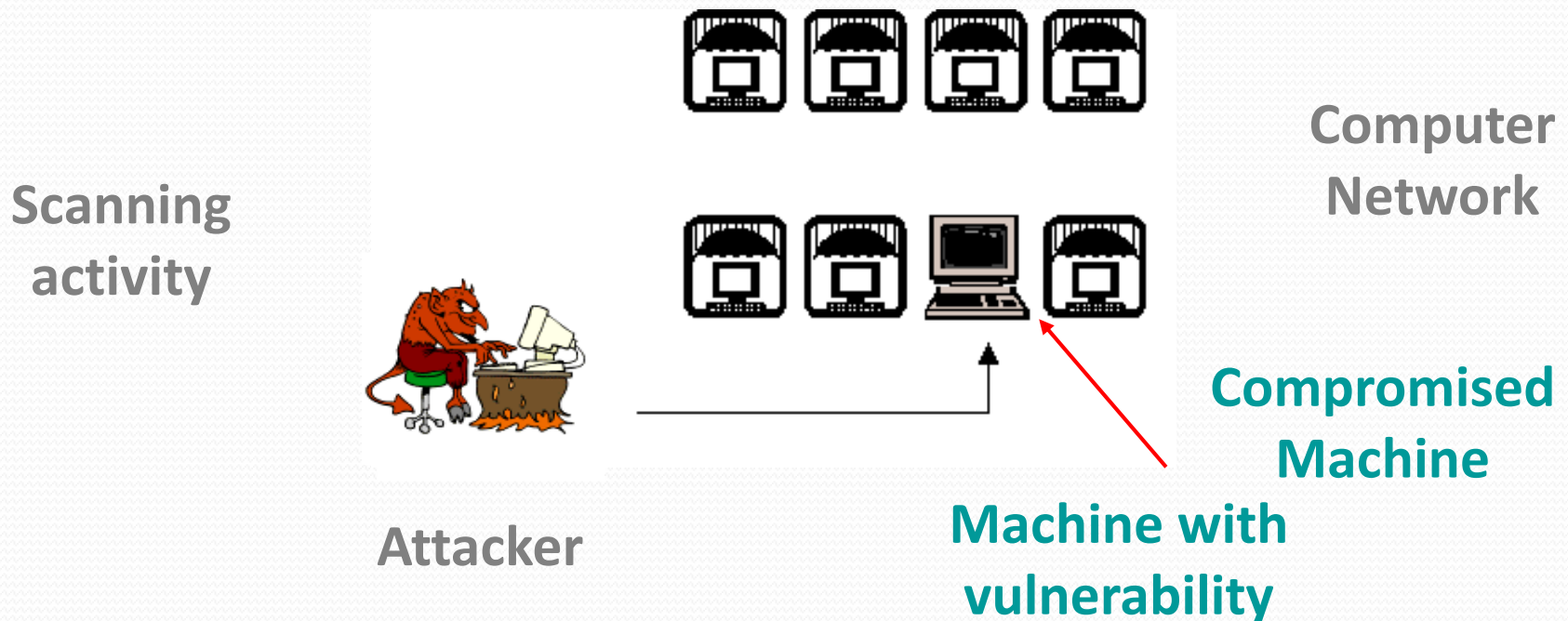
Intrusion Detection System(IDS)

- combination of software and hardware that attempts to perform intrusion detection
- raise the alarm when possible intrusion or suspicious patterns are observed

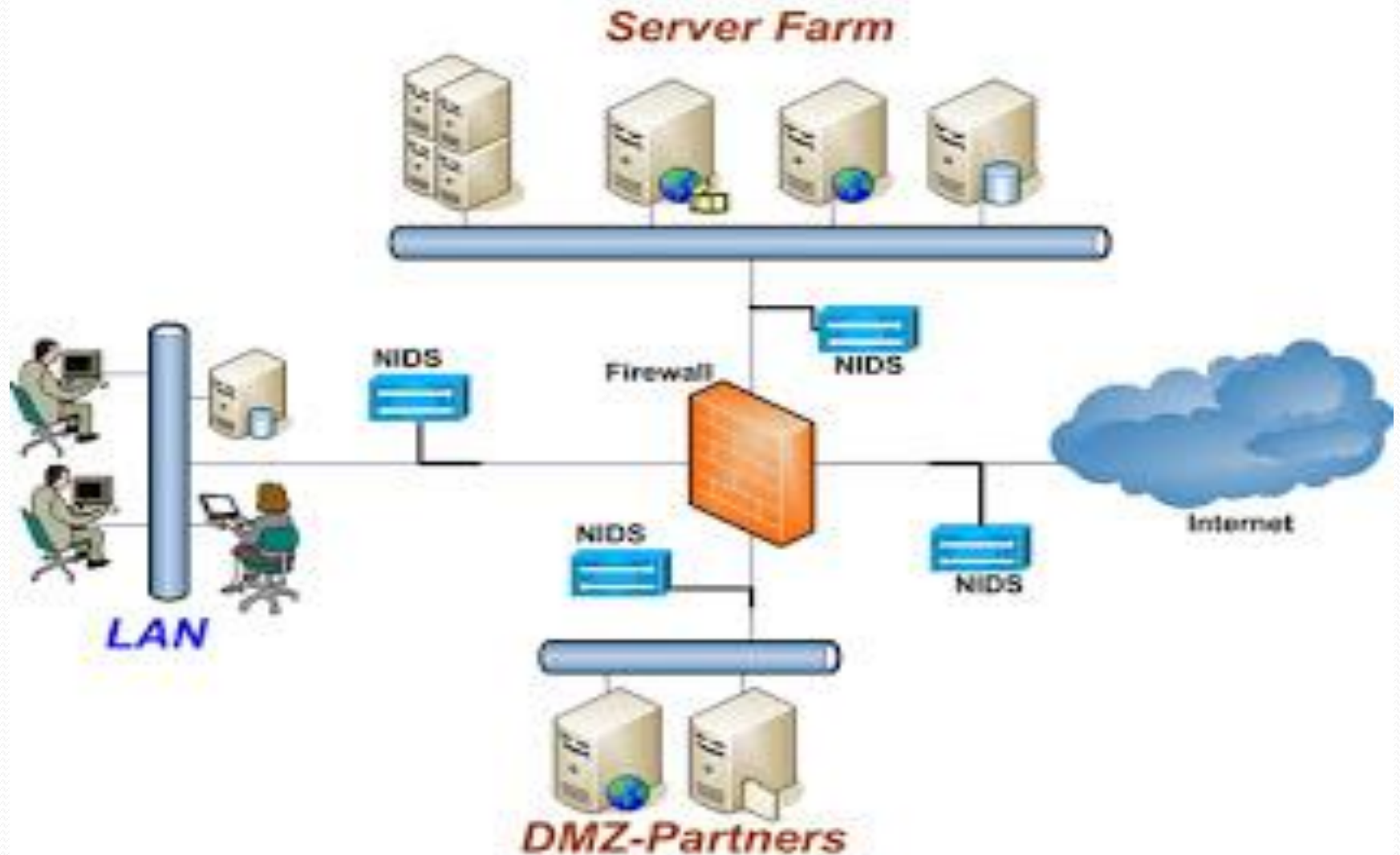


What are Intrusions?

- ◆ Intrusions are actions that attempt to bypass security mechanisms of computer systems. They are usually caused by:
 - Attackers accessing the system from Internet
 - Insider attackers - authorized users attempting to gain and misuse non-authorized privileges
- ◆ Typical intrusion scenario

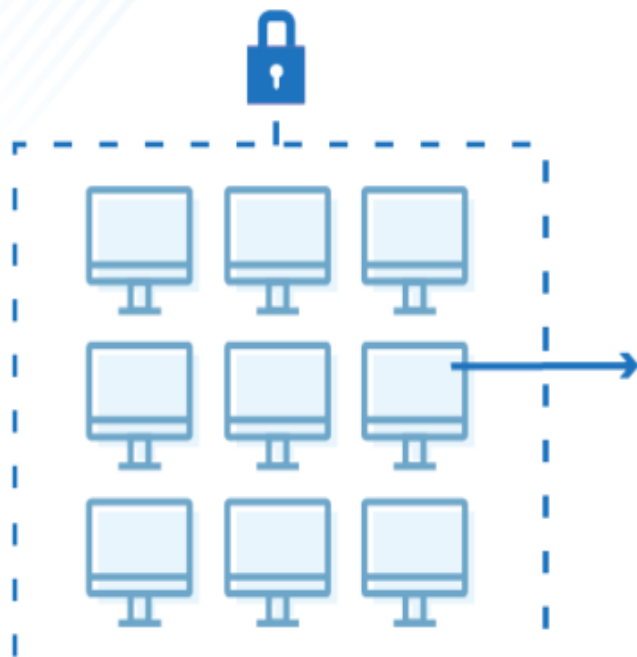


Network Diagram of Enterprise Network



Types of IDS

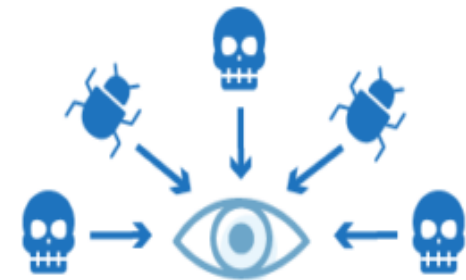
What Does an Intrusion Detection System Do?



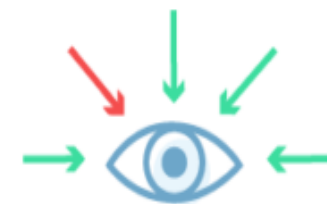
Network Intrusion
Detection



Host Intrusion
Detection



Signature-based
Detection



Anomaly-Based
Detection

Data Mining for Intrusion Detection

- *Signature-based / Misuse detection*

- ◆ Building predictive models from labeled data sets (instances are labeled as “normal” or “intrusive”) to identify known intrusions
- ◆ High accuracy in detecting many kinds of known attacks
- ◆ Cannot detect unknown and emerging attacks

- *Anomaly detection*

- ◆ Detect novel attacks as deviations from “normal” behavior
- ◆ Potential high false alarm rate - previously unseen (yet legitimate) system behaviors may also be recognized as anomalies

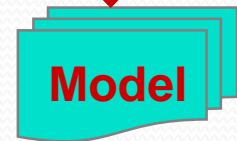
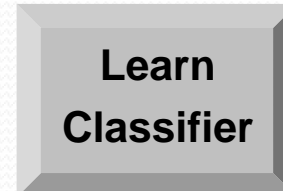
Data Mining for Signature-based IDS

Misuse Detection – Building Predictive Models

categorical
temporal
categorical
continuous
class

Tid	SrcIP	Start time	Dest IP	Dest Port	Number of bytes	Attack
1	206.135.38.95	11:07:20	160.94.179.223	139	192	No
2	206.163.37.95	11:13:56	160.94.179.219	139	195	No
3	206.163.37.95	11:14:29	160.94.179.217	139	180	No
4	206.163.37.95	11:14:30	160.94.179.255	139	199	No
5	206.163.37.95	11:14:32	160.94.179.254	139	19	Yes
6	206.163.37.95	11:14:35	160.94.179.253	139	177	No
7	206.163.37.95	11:14:36	160.94.179.252	139	172	No
8	206.163.37.95	11:14:38	160.94.179.251	139	285	Yes
9	206.163.37.95	11:14:41	160.94.179.250	139	195	No
10	206.163.37.95	11:14:44	160.94.179.249	139	163	Yes

Tid	SrcIP	Start time	Dest IP	Number of bytes	Attack
1	206.163.37.81	11:17:51	160.94.179.208	150	No
2	206.163.37.99	11:18:10	160.94.179.235	208	No
3	206.163.37.55	11:34:35	160.94.179.221	195	Yes
4	206.163.37.37	11:41:37	160.94.179.253	199	No
5	206.163.37.41	11:55:19	160.94.179.244	181	Yes

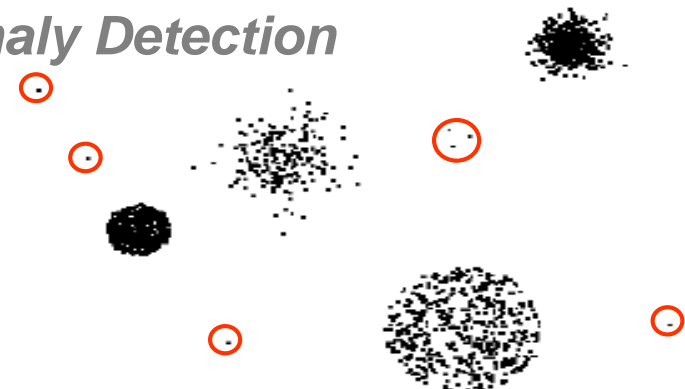


Summarization of attacks using association rules

Rules Discovered:

{Src IP = 206.163.37.95, Dest Port = 139, Bytes ∈ [150, 200]} --> {ATTACK}

Anomaly Detection



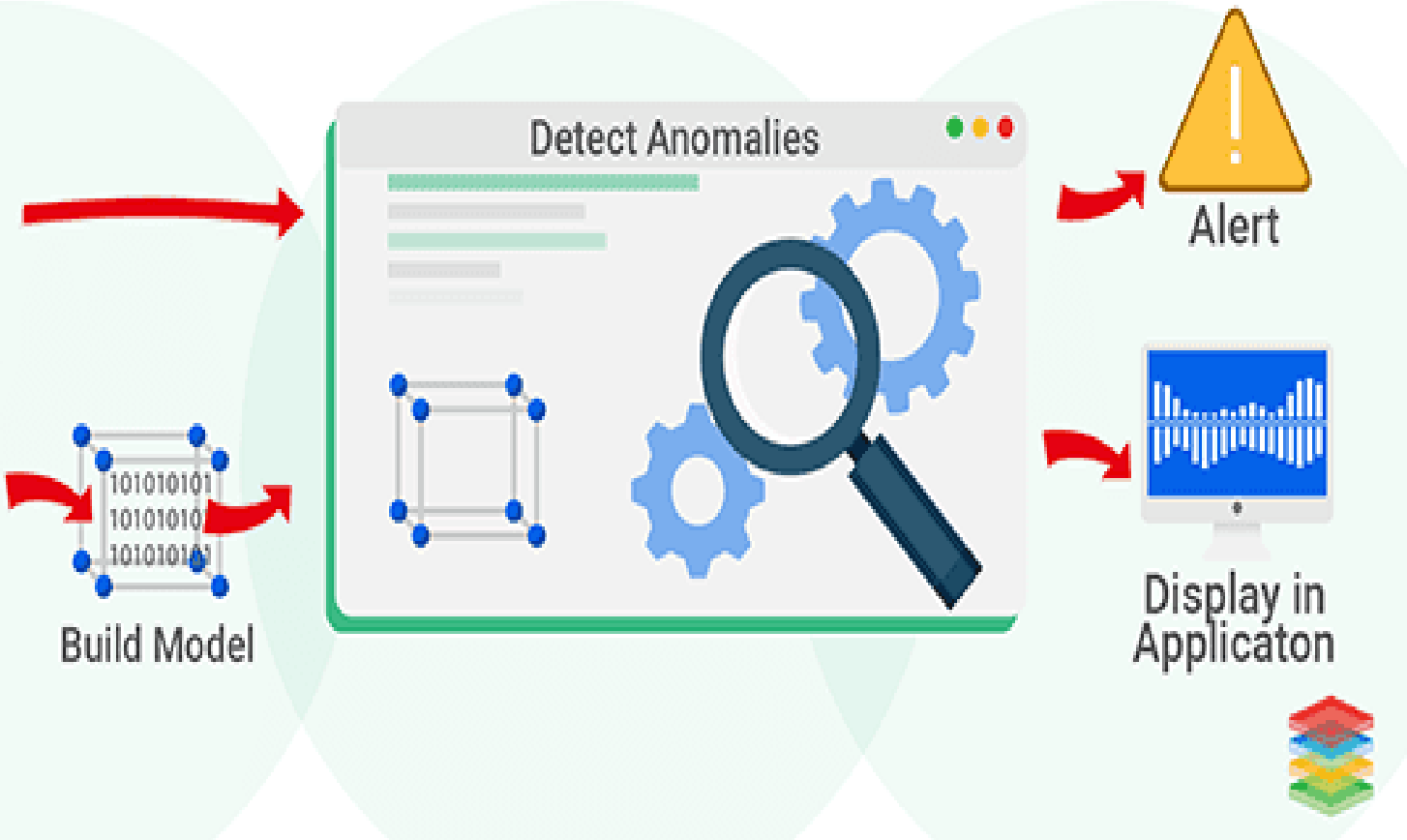
Real Time Anomaly Detection

Live Data

101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101

Historical Data

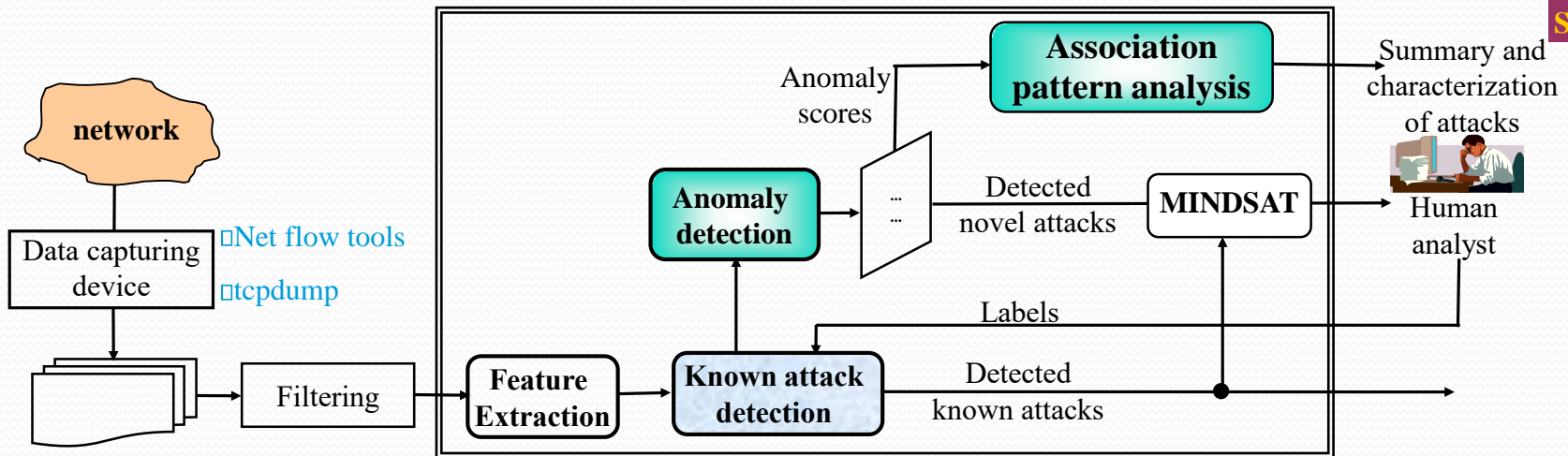
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101



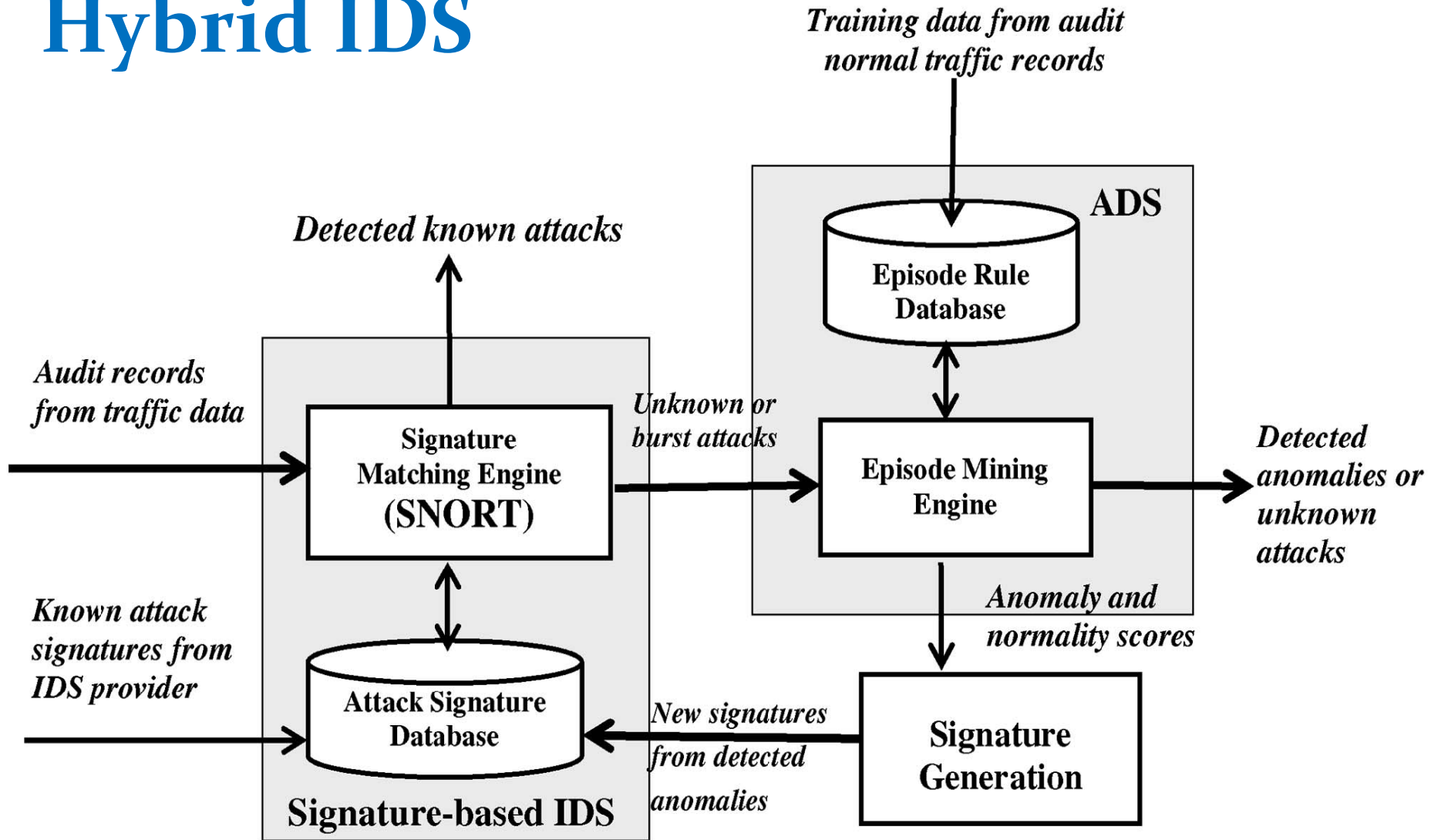
Hybrid IDS: signature-based + anomaly detection

- Anomaly detection was used at U of Minnesota and Army Research Lab to detect various intrusive/suspicious activities
- Many of these could not be detected using widely used intrusion detection tools like SNORT
- Anomalies/attacks picked by *MINDS*
 - Scanning activities
 - Non-standard behavior
 - Policy violations
 - Worms

MINDS – Minnesota Intrusion Detection System



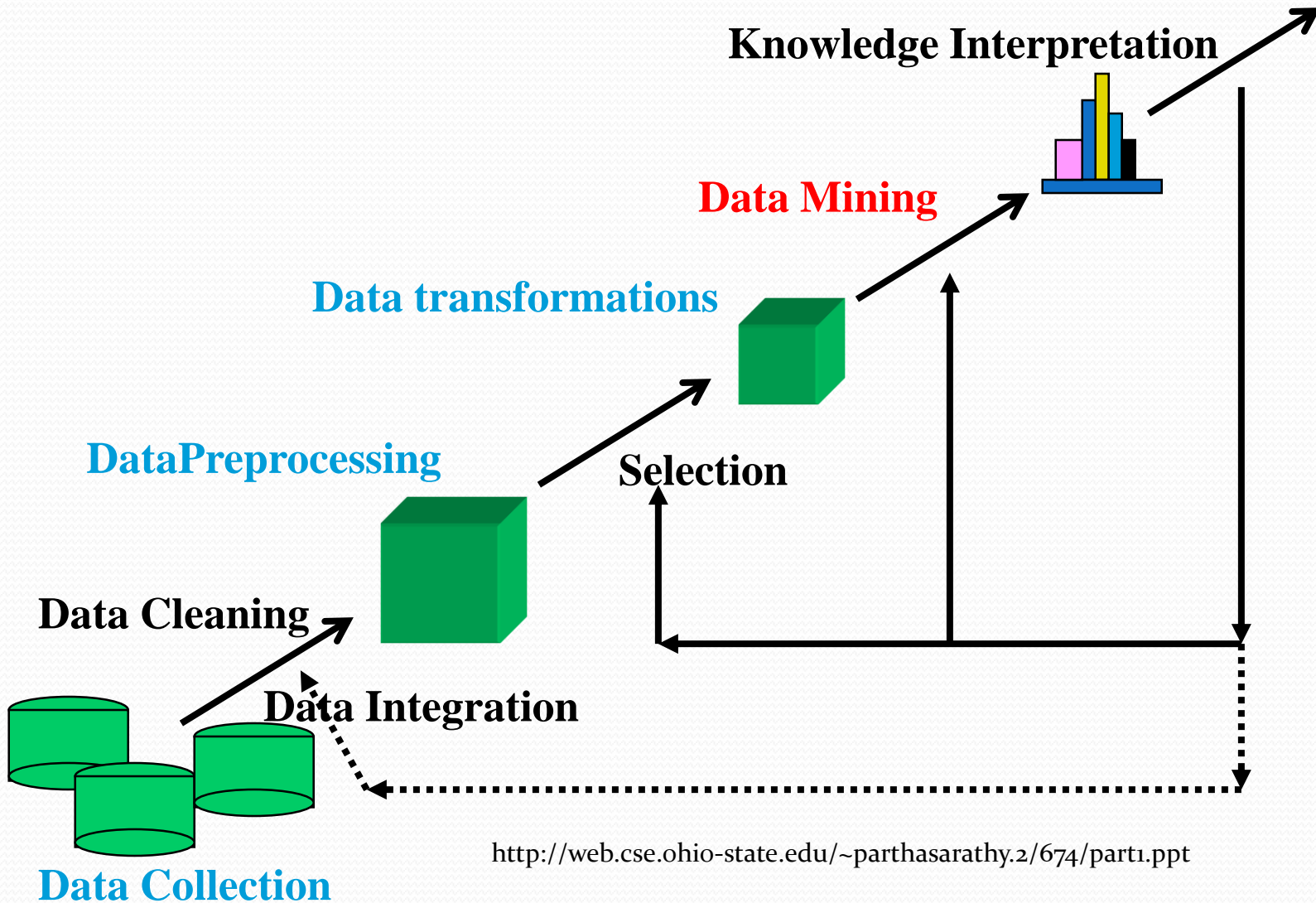
Hybrid IDS



Data Mining

Data mining:
the core of knowledge discovery process.

Knowledge

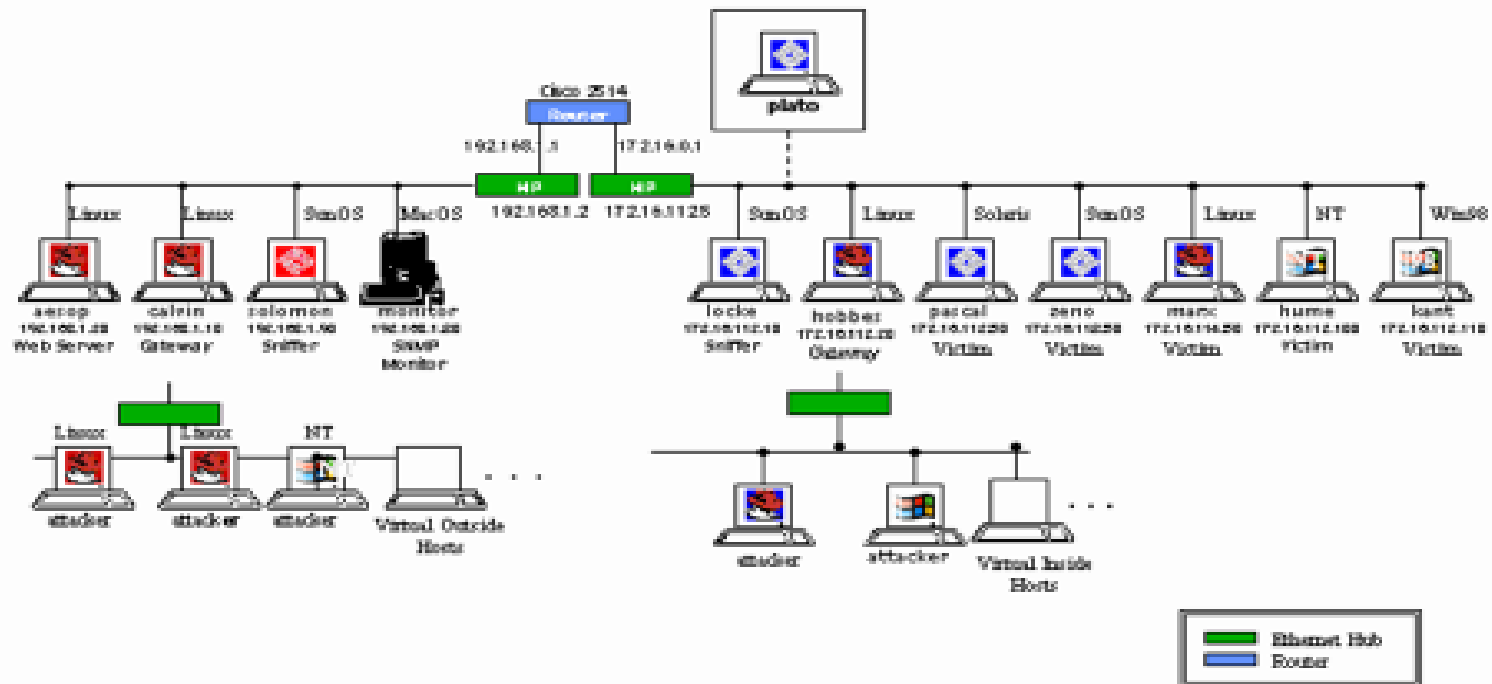


Step 1. Intrusion Data Collection

(raw data : real network traffics)



Simulation Network 99



Step 2. Data Pre-Processing

DARPA-MIT dataset : real time network traffics

- DARPA 1998 data set
 - Simulated nine weeks of raw TCP dump data
 - Probing attacks, DoS attacks, U2R, R2L attacks

```
user@host:~$ sudo tcpdump --interface=ens33 -n host 192.168.111.1
tcpdump: verbose output suppressed, use -v or -vv for full protocol decode
listening on ens33, link-type EN10MB (Ethernet), capture size 262144 bytes
23:55:40.546464 IP 192.168.111.1 > 192.168.111.209: ICMP echo request, id 64017, seq 0, length 64
23:55:40.546517 IP 192.168.111.209 > 192.168.111.1: ICMP echo reply, id 64017, seq 0, length 64
23:55:41.551452 IP 192.168.111.1 > 192.168.111.209: ICMP echo request, id 64017, seq 1, length 64
23:55:41.551485 IP 192.168.111.209 > 192.168.111.1: ICMP echo reply, id 64017, seq 1, length 64
23:55:42.556206 IP 192.168.111.1 > 192.168.111.209: ICMP echo request, id 64017, seq 2, length 64
23:55:42.556243 IP 192.168.111.209 > 192.168.111.1: ICMP echo reply, id 64017, seq 2, length 64
23:55:43.558055 IP 192.168.111.1 > 192.168.111.209: ICMP echo request, id 64017, seq 3, length 64
23:55:43.558094 IP 192.168.111.209 > 192.168.111.1: ICMP echo reply, id 64017, seq 3, length 64
23:55:43.955857 IP 192.168.111.1.53861 > 192.168.111.209.80: Flags [SEW], seq 3194582235, win 65535, options [mss 14
60,nop,wscale 6,nop,nop,TS val 243647685 ecr 0,sackOK,eol], length 0
23:55:43.955909 IP 192.168.111.209.80 > 192.168.111.1.53861: Flags [S.E], seq 4099266365, ack 3194582236, win 65160,
options [mss 1460,sackOK,TS val 1285093713 ecr 243647685,nop,wscale 7], length 0
23:55:43.956230 IP 192.168.111.1.53861 > 192.168.111.209.80: Flags [.), ack 1, win 2058, options [nop,nop,TS val 243
647685 ecr 1285093713], length 0
23:55:43.956250 IP 192.168.111.1.53861 > 192.168.111.209.80: Flags [P.), seq 1:80, ack 1, win 2058, options [nop,nop
,TS val 243647685 ecr 1285093713], length 79: HTTP: GET / HTTP/1.1
23:55:43.956385 IP 192.168.111.209.80 > 192.168.111.1.53861: Flags [.), ack 80, win 509, options [nop,nop,TS val 128
5093713 ecr 243647685], length 0
```

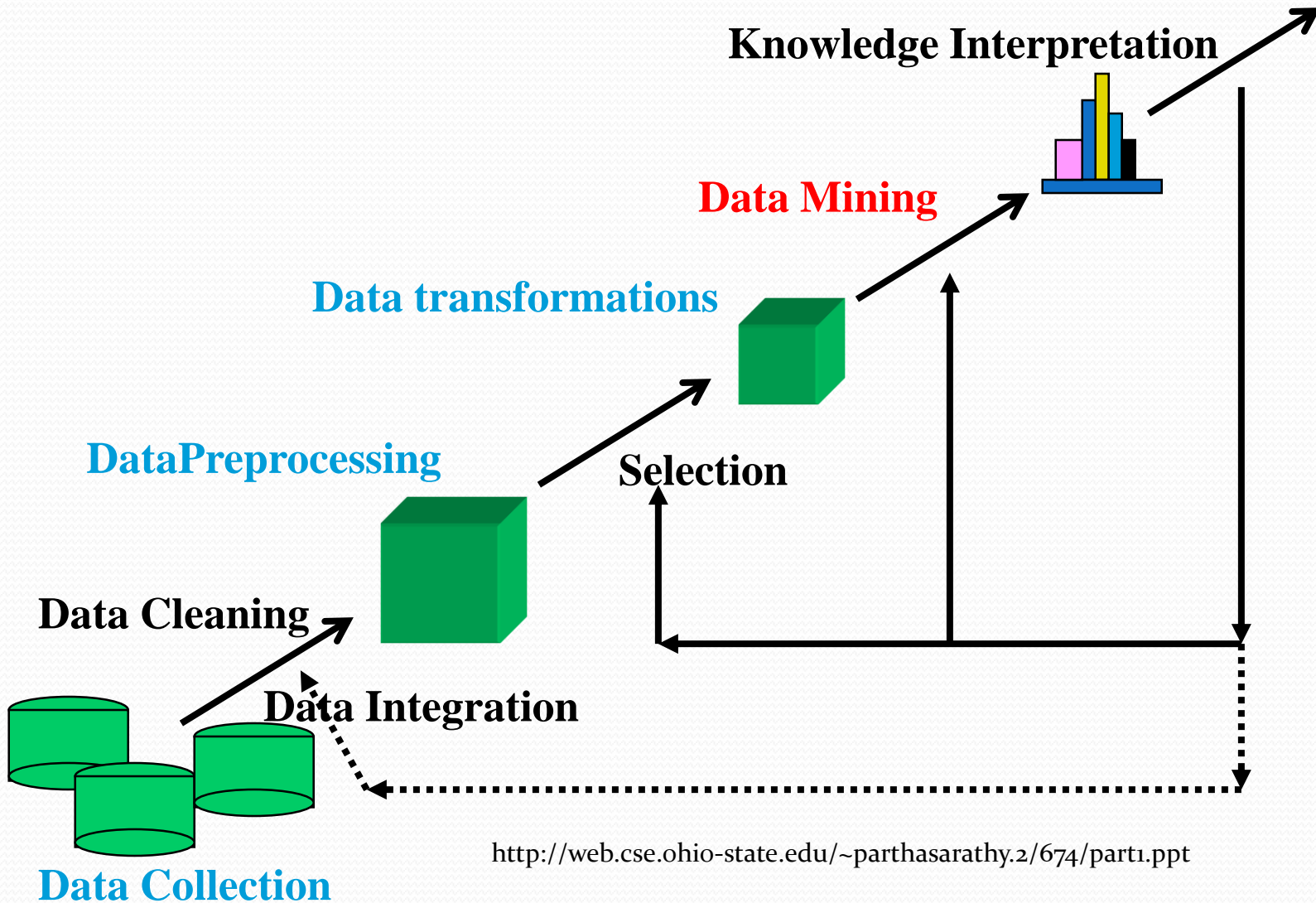
Intrusion Datasets (ready to used)

- Darpa-Intrusion Dataset 1998
- KDD Cup Intrusion Data 1999
- NSL-KDD Intrusion Dataset
- Kyoto Intrusion Dataset 2006
- Intrusion Detection Evaluation Dataset (CICIDS2017)
 - Android Botnet 2015
 - Android Adware 2017
 - Botnet 2014
 - Denial of Service Attack 2017
 - Distributed Denial of Service Attack 2017 & 2019
 - IDS 2019
 - DNS over HTTPS attacks 2020
- <https://www.unb.ca/cic/datasets/index.html>

Data Mining

Data mining:
the core of knowledge discovery process.

Knowledge



Step 3 :Data Transformation

- Three groups of features

- Basic features of individual TCP connections

- source & destination IP *Features 1 & 2*
- source & destination port *Features 3 & 4*
- Protocol *Feature 5*
- Duration *Feature 6*
- Bytes per packets *Feature 7*
- number of bytes *Feature 8*

<i>dst ...</i>	<i>service ...</i>	<i>flag</i>		<i>dst ...</i>	<i>service ...</i>	<i>flag</i>	<i>%S0</i>
h1	http	S0	syn flood	h1	http	S0	70
h1	http	S0		h1	http	S0	72
h1	http	S0		h1	http	S0	75
h2	http	S0	normal	h2	http	S0	0
h4	http	S0		h4	http	S0	0
h2	ftp	S0		h2	ftp	S0	0

existing features
useless

construct features with
high information gain

- Time based features

- For the same source (*destination*) IP address, number of unique destination (*source*) IP addresses inside the network *in last T seconds* – *Features 9 (13)*
- Number of connections from source (*destination*) IP to the same destination (*source*) port *in last T seconds* – *Features 11 (15)*

- Connection based features

- For the same source (*destination*) IP address, number of unique destination (*source*) IP addresses inside the network *in last N connections* - *Features 10 (14)*
- Number of connections from source (*destination*) IP to the same destination (*source*) port *in last N connections* - *Features 12 (16)*

Step 3 :Data Transformation

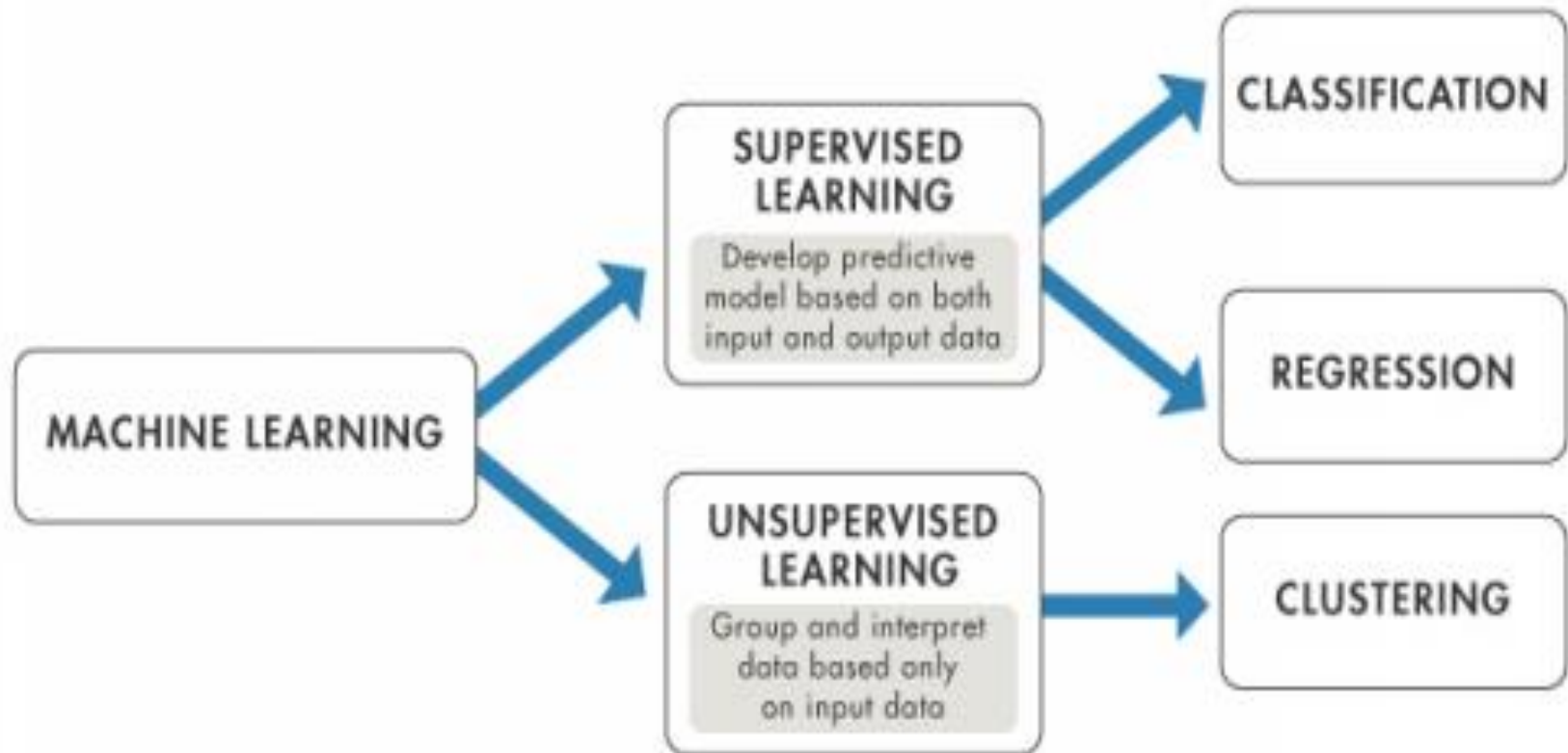
KDD Cup 99 Data Set : 41 attributes + label

S.NO	FEATURE NAME	S.NO	FEATURE NAME	Category	Class label(attack) in dataset
1	Duration	22	Is_guest_login		
2	Protocol type	23	Count		
3	Service	24	Serror_rate		
4	Src_byte	25	Rerror_rate		
5	Dst_byte	26	Same_srv_rate		
6	Flag	27	Diff_srv_rate	DOS-Denial of service	back,land, pod, neptune, smurf, teardrop
7	Land	28	Srv_count		
8	Wrong_fragment	29	Srv_serror_rate		
9	Urgent	30	Srv_rerror_rate	R2L-Remote to local	ftp_write, guess_passwd, imap, multihop, phf, spy, warezclient
10	Hot	31	Srv_diff_host_u		
11	Num_failed_logins	32	Dst_host_coun	U2R-User to root	Buffer_overflow, loadmodule, perl, rootkit
12	Logged_in	33	Dst_host_srv_c		
13	Num_compromised	34	Dst_host_same		
14	Root_shell	35	Dst_host_diff_s		
15	Su_attempted	36	Dst_host_same	Probe	Ipsweep, nmap, portsweep, satan
16	Num_root	37	Dst_host_srv_ditt_host_rate		
17	Num_file_creations	38	Dst_host_serror_rate		
18	Num_shells	39	Dst_host_srv_serror_rate		
19	Num_access_shells	40	Dst_host_rerror_rate		
20	Num_outbound_cmds	41	Dst_host_srv_rerror_rate		
21	Is_hot_login				

Data Mining Techniques



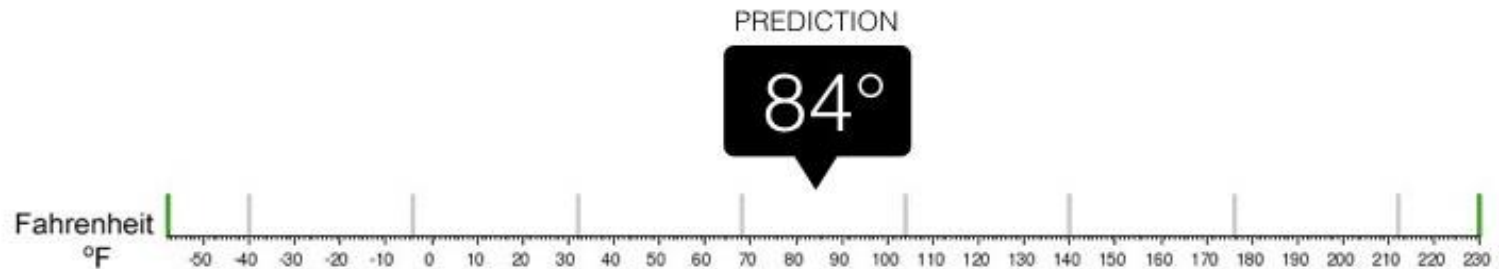
Supervised Learning vs Unsupervised Learning





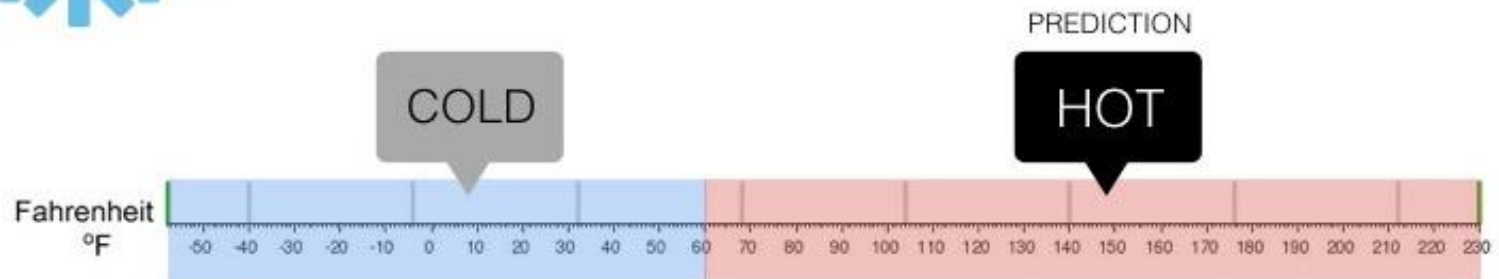
Regression

What is the temperature going to be tomorrow?



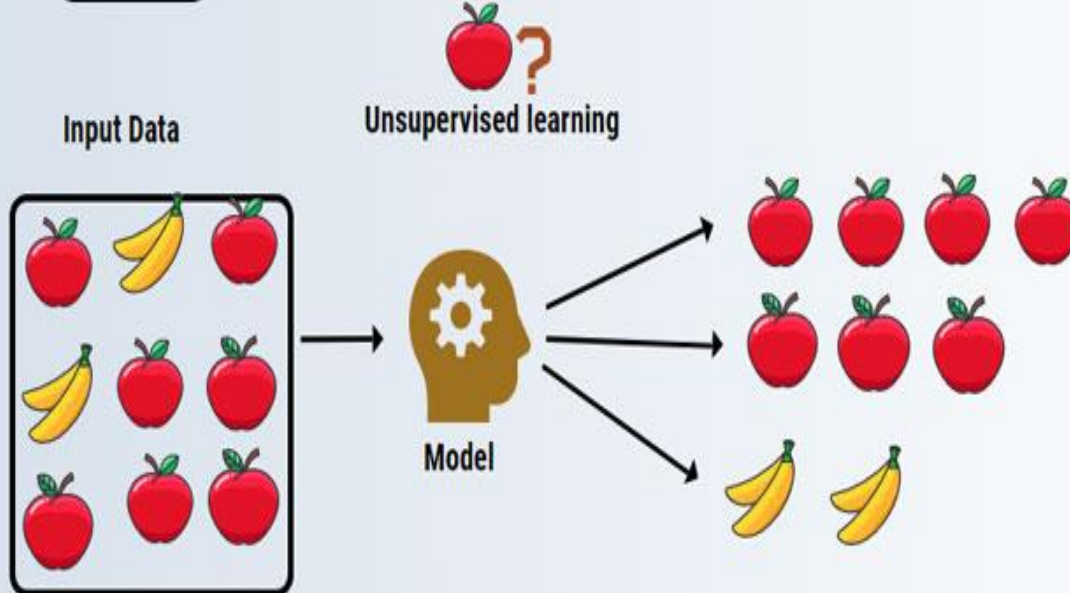
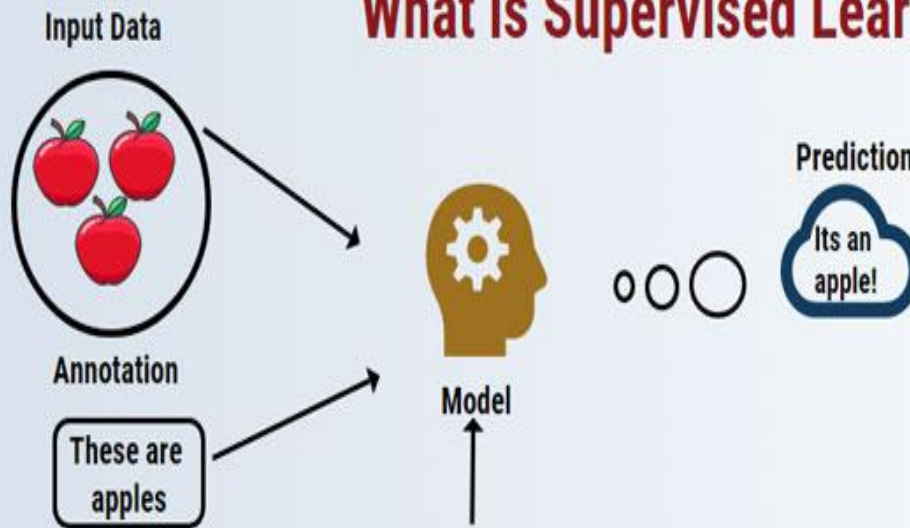
Classification

Will it be Cold or Hot tomorrow?



Supervised Learning

What is Supervised Learning?



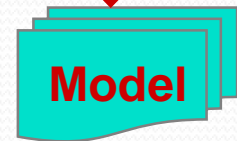
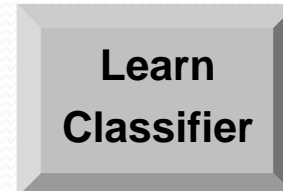
Data Mining for Signature-based IDS

Misuse Detection – Building Predictive Models

categorical
temporal
categorical
continuous
class

Tid	SrcIP	Start time	Dest IP	Dest Port	Number of bytes	Attack
1	206.135.38.95	11:07:20	160.94.179.223	139	192	No
2	206.163.37.95	11:13:56	160.94.179.219	139	195	No
3	206.163.37.95	11:14:29	160.94.179.217	139	180	No
4	206.163.37.95	11:14:30	160.94.179.255	139	199	No
5	206.163.37.95	11:14:32	160.94.179.254	139	19	Yes
6	206.163.37.95	11:14:35	160.94.179.253	139	177	No
7	206.163.37.95	11:14:36	160.94.179.252	139	172	No
8	206.163.37.95	11:14:38	160.94.179.251	139	285	Yes
9	206.163.37.95	11:14:41	160.94.179.250	139	195	No
10	206.163.37.95	11:14:44	160.94.179.249	139	163	Yes

Tid	SrcIP	Start time	Dest IP	Number of bytes	Attack
1	206.163.37.81	11:17:51	160.94.179.208	150	No
2	206.163.37.99	11:18:10	160.94.179.235	208	No
3	206.163.37.55	11:34:35	160.94.179.221	195	Yes
4	206.163.37.37	11:41:37	160.94.179.253	199	No
5	206.163.37.41	11:55:19	160.94.179.244	181	Yes

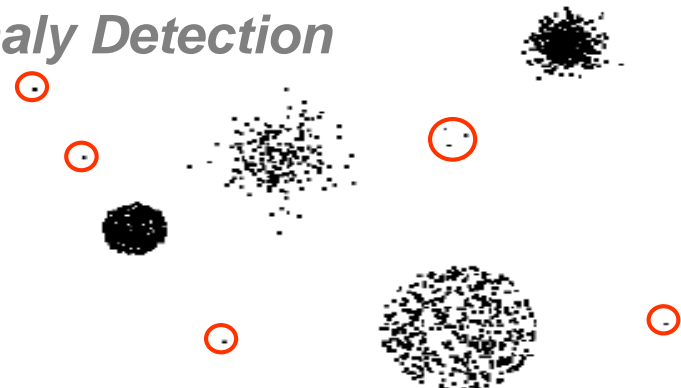


Summarization of attacks using association rules

Rules Discovered:

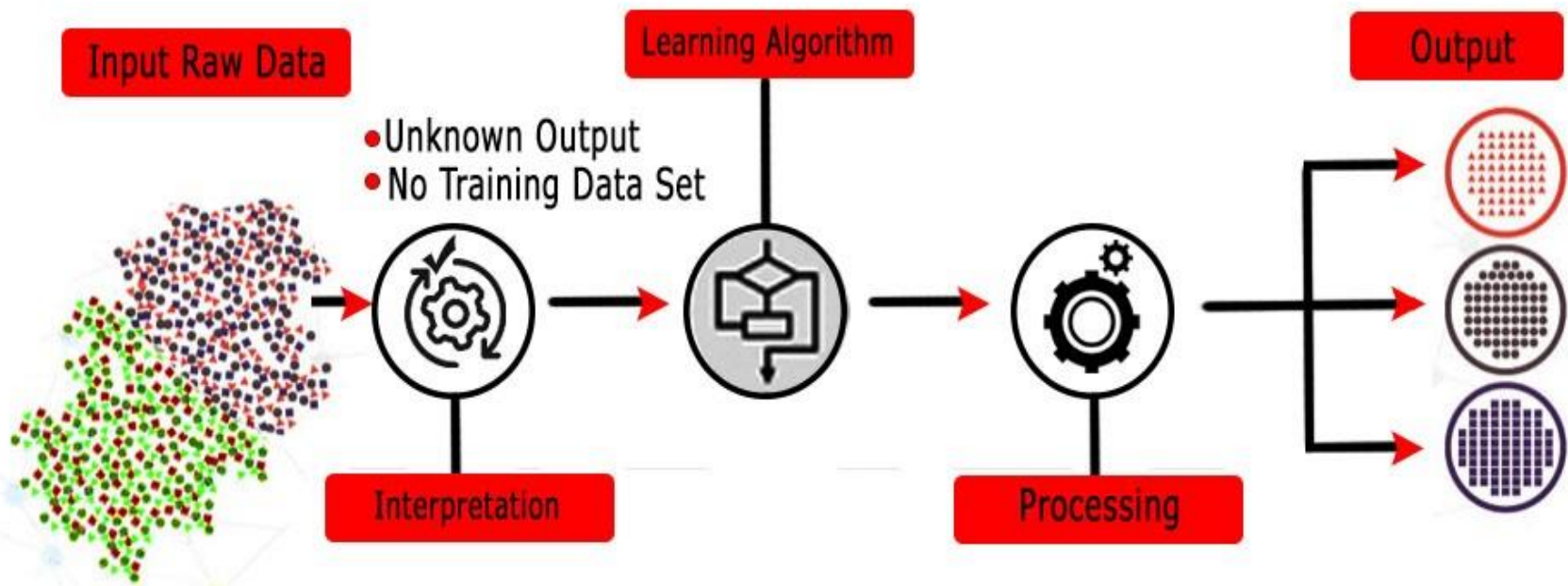
{Src IP = 206.163.37.95, Dest Port = 139, Bytes ∈ [150, 200]} --> {ATTACK}

Anomaly Detection



Unsupervised Learning

Unsupervised Learning



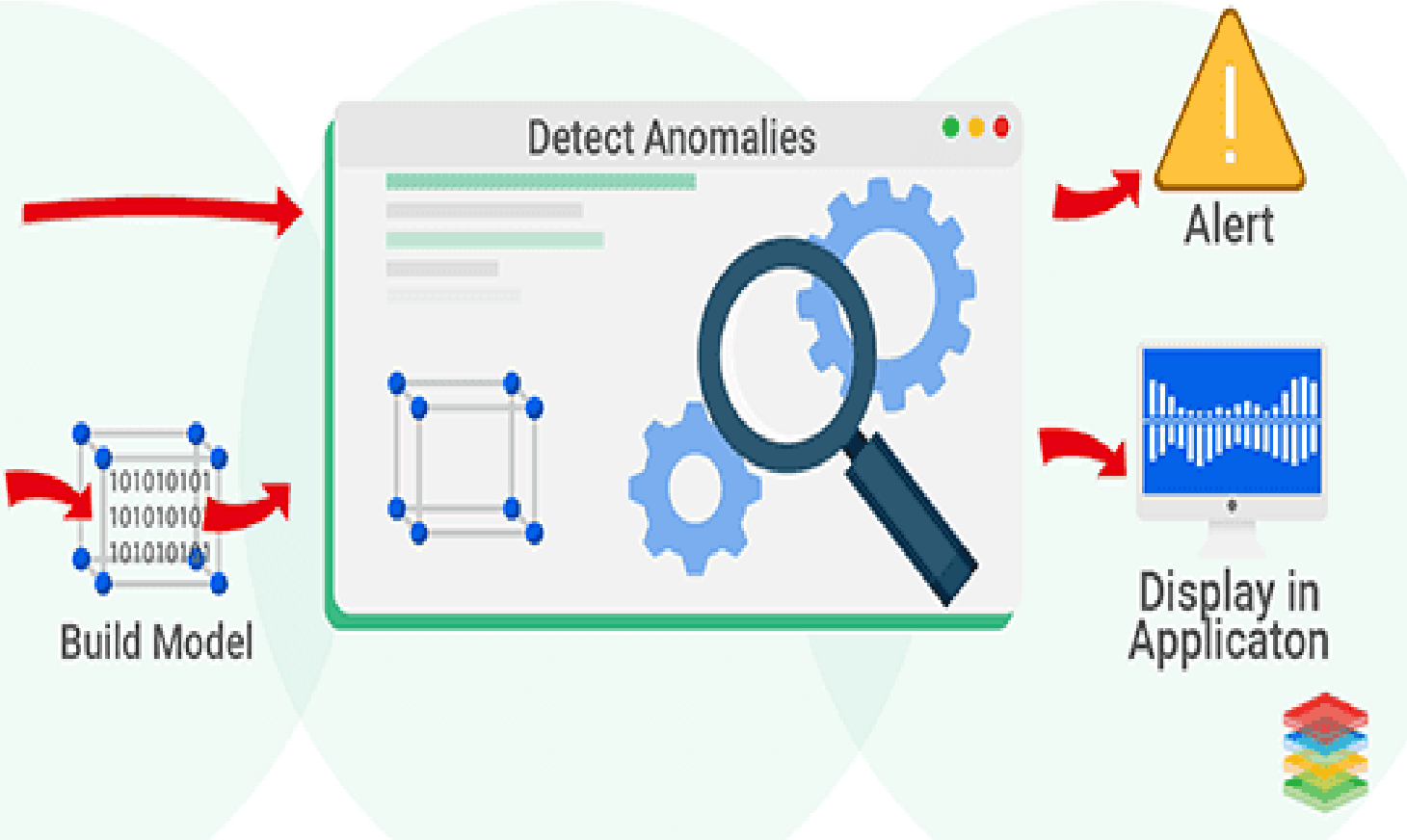
Real Time Anomaly Detection

Live Data

```
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
```

Historical Data

```
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
101010101010101
```

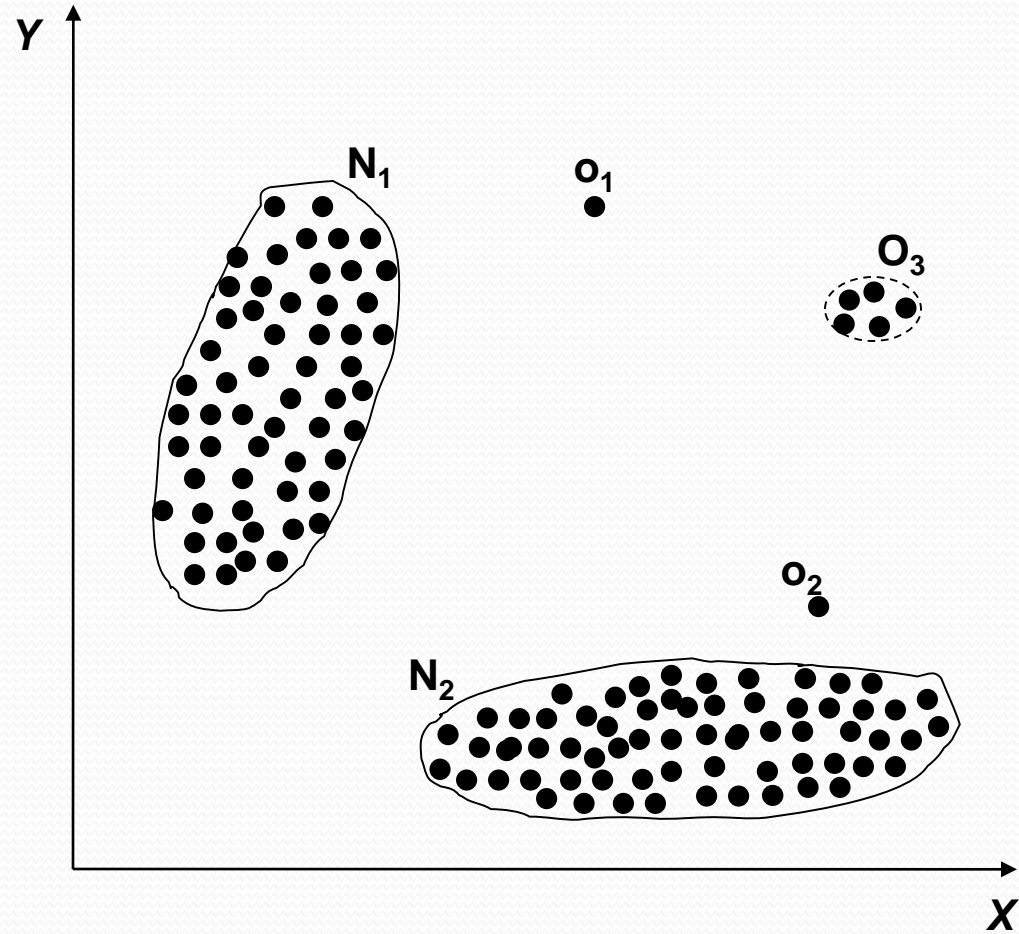


What are Anomalies?

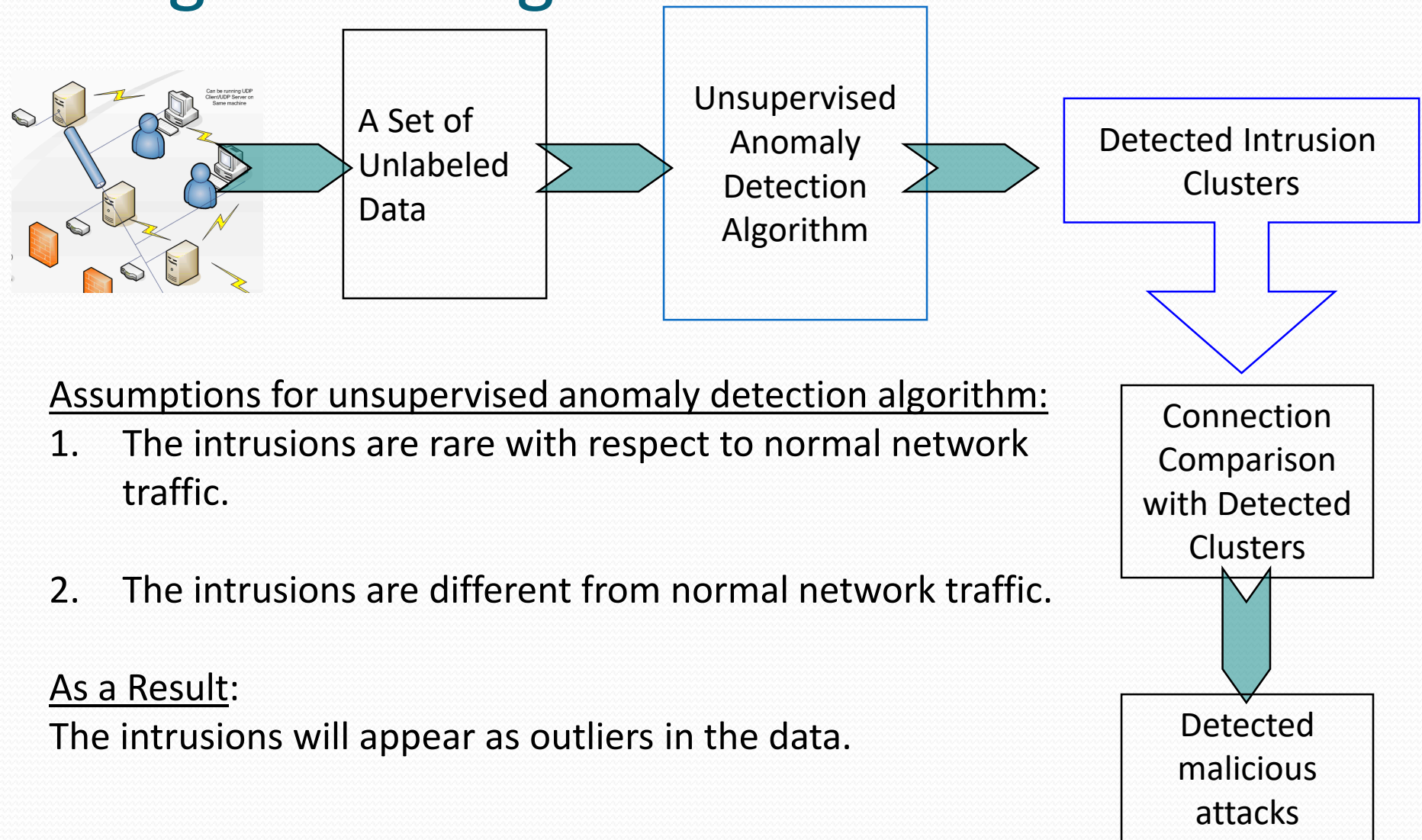
- Anomaly is a pattern in the data that does not conform to the expected behavior
- Also referred to as outliers, exceptions, peculiarities, surprise, etc.
- Anomalies translate to significant (often critical) real life entities
 - Cyber intrusions
 - Credit card fraud

Simple Example

- N_1 and N_2 are regions of normal behavior
- Points o_1 and o_2 are anomalies
- Points in region O_3 are anomalies



Using Clustering for Intrusion Detection



Assumptions for unsupervised anomaly detection algorithm:

1. The intrusions are rare with respect to normal network traffic.
2. The intrusions are different from normal network traffic.

As a Result:

The intrusions will appear as outliers in the data.

Using Clustering for Intrusion Detection

Once data is clustered, all of the instances that appear in small clusters are labeled as anomalies because;

- The normal instances should form large clusters compared to the intrusions,
- Malicious intrusions and normal instances are qualitatively different, so they do not fall into the same cluster.

Intrusion cluster

