

CS 358 ESE ASSIGNMENT

Navyasree
1801CS06
CS358
16

① Match-table entries:-

$S_1 \rightarrow$ match	Action
IP Src = 128.119.16	Forward (4)
Ip protocol = TCP	
IP dst = 128.121.16	Forward (2)
$S_2 \rightarrow$ match	Action
Ingress port = 3	
IP Src = 128.119.16	Forward (2)
IP dst = 128.121.16	
IP protocol = UDP	
$S_3 \rightarrow$ match	Action
Ingress port = 2	Forward (1)
IP dst = 128.121.16	
Ingress port = 4	Forward (1)
IP dst = 128.121.16	
$S_4 \rightarrow$ match	Action
Ingress port = 2	
IP Src = 128.119.16	Forward (3)
IP dst = 128.121.16	
IP protocol = TCP	

- (a) For router S₁, Dst port can be "any"
Unless router needs to block packet from exiting based on context, dst port will be 'any'
- (b) For router S₁, Both TCP & UDP are used.
for interface 2, UDP used for path: S₁ → S₂ → S₃
For interface 4, TCP used for path: S₁ → S₄ → S₃
- (c) For router S₂, IP src will be "128.119/16"
Source will be subnet of where packet originated
- (d) Here, S₃ forwards packets to 128.121/6 through interface 1. So Action is 'forward'
- (e) For router S₄, packets needs to be forwarded to interface "3" for it to reach S₃.

(2) AF

(a) The distance vector from router when algorithm converges, 'x', is

" $\{11, 8, 9, 10, 5\}$ " (By finding Shortest distance from 'x')

(b) Initial distance vectors for router 'y' is

$$(u, v, w, x, y) = (\infty, \infty, \infty, 5, 0)$$

(c) Name of problem which occurs when link costs increase is "Count to infinity problem"

(3) AF

(a) we know, Shortest distance for 'x' to 'v' is 3
So, for link z, the cost associated with the link is

$$3 \Rightarrow \underline{\underline{x = 3}}$$

(b) we know, Shortest path from x to z = 5.
which is along path $x \rightarrow w \rightarrow z$

There is no possibility we find link 'y'

Hence, n/a

(A)

Given,

Two self learning switches (switches in which they learn which hosts can be through which interface)

Through table entries for both switches,

(a) At $t=4$, Nodes communicating are C

In given table 1, there is only packet transmission record to C at TTL = 4

(b) At $t=8$, Nodes communicating are C, F

for TTL = 8, packet transmission between C & F as recorded in table 2.

(c) At $t=3$, Nodes communicating are J

packet transmission through J is recorded both in table 1 & table 2 at packet whose TTL = 3

(d) At $t=2$, Nodes communicating are H

Packet transmission through H is recorded in both table 1 & table 2 whose TTL = 2

(5) A:

At $t=1$: $F \rightarrow I$

So packet transmissions recorded after this step,

MAC interface	TTL	MAC interface	TTL
F	8	I	1
I	7	F	1

At $t=2$: $C \rightarrow J$

After $C \rightarrow J$, table entries look like

MAC interface	TTL	MAC interface	TTL
F	8	I	1
I	12	J	2
C	8	C	3
J	13	J	2

At $t=3$: $E \rightarrow G$

After $E \rightarrow G$, table entries look like.

MAC interface	TTL	MAC interface	TTL
F	8	I	1
I	12	J	2
C	8	C	3
J	13	J	2
E	8	E	3
G	10	G	3

At $t=4$: $F \rightarrow C$

After these transmission from $F \rightarrow C$

Final table looks,

MAC Interface	TTL	MAC Interface	TTL
F	8	F	6
I	12	I	7
C	8	C	3
J	13	J	7
E	8	E	5
G	10	G	7
F	8	F	6
C	8	C	3

(a) At $t=1$, we have seen in step 1 Source entry

$(F, 6)$ is added as ^{Source} entry to Switch 1

(b) At $t=1$, in step 1 we have seen that

$(I, 7)$ is added as destination entry to Switch 1

(c) At $t=2$, we have seen in step 2

$(J, 13)$ is destination entry to Switch 2

(d) At $t=3$, we have Seen in Step 3

$(E_1, 5)$ is Source entry for switch 1

(e) At $t=3$, we have Seen in Step 3

$(G_1, 10)$ is destination entry for Switch 2

⑥ A:

Cyclic Redundancy check algorithm (CRC):

CRC is used to detect accidental changes to data transmitted via networks. CRC involves binary division of data bits being sent by a pre-determined divisor agreed by communicating system. Divisor is generated using co-efficient of polynomial. We append 'r' bits to datapayload(D) and store the remainder of (D) when divided by G and send the data (D+remainder).

Now,

We divide the data received on other side with the same divisor (G)

If remainder is '0' \Rightarrow No error in transmission

If remainder \neq '0' \Rightarrow Error in transmission

Given,

$$r = 3$$

Generator / Divisor (G) = 1001

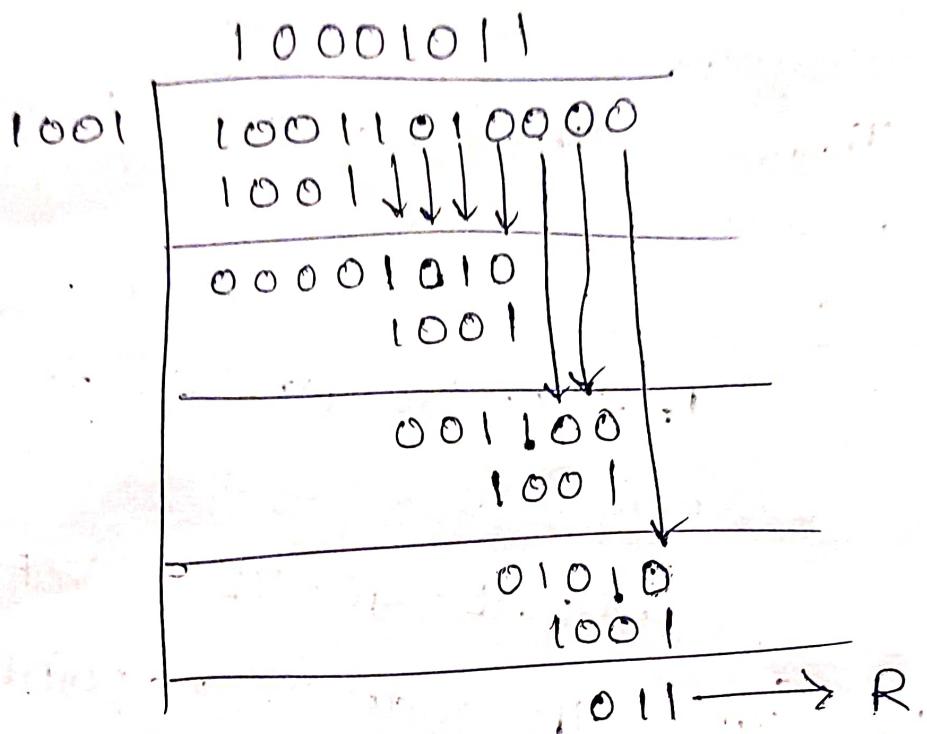
Datapayload (D) = 100110010

We add 'r' zeros to D, divide by G,

$$D + R = 10011010000 \text{ ('+' means appended)}$$

$$R = 1001$$

On division,



We got that, Remainder (R) = 011

We will add 111 to D & Send the data

∴ CRC bits (R) associated with Datapayload (D) = 011

Q7 At

Given,

Source and destination addresses of IP datagram & Ethernet address (MAC address)
Ethernet is physical address, IP is logical address.

(a) Source mac address at point 5:

1A-79-BC-E9-E6-3A

(b) Destination mac address at point 5:

DO-AE-40-72-CC-FO

(c) Source IP address at point 5:

128.119.58.43

(d) Destination IP address at point 5:

128.119.190.195

(e) Source mac address at point 2:

E6-E3-64-04-2A-03

(f) Destination mac address at point 2:

C2-32-12-7D-4E-86

These are required mac & IP addresses
at point 5/2

Ques:

Given,

payload has 10 8-bit values.

Method:

calculate number of one's from top to bottom
& left to right if the count is odd mark
it as 1 else mark it as 0

1	0	1	0	0	1	1	1	0	1	0	0	1	0	0	0	1
1	0	1	1	1	1	0	1	0	0	1	1	0	0	0	0	0
1	0	1	1	1	0	1	0	0	1	0	1	0	1	1	0	1
0	1	1	0	0	0	0	1	1	0	0	1	1	0	0	0	0
0	0	1	0	1	1	0	1	1	1	0	1	0	0	0	0	0
⇒	1	1	1	0	1	1	0	0	0	1	0	1	1	1	1	0

(a) On Counting from top to bottom,

2-Dimensional parity for 16 columns $\rightarrow 1110110001011110$

(b) On Counting from left to right,

2 - Dimensional Parity for 5 rows $\rightarrow 10100$

Q1:

Let us now look at the Oscillation problem faced in traditional Computer networks.

We will consider example Dijkstra's algorithm where oscillation occurs.

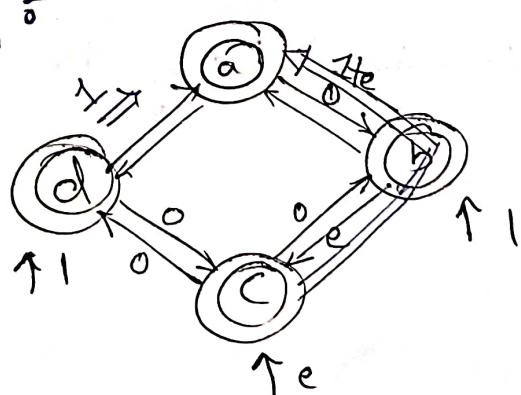
Link costs between any two nodes in the network are dependent on amount of traffic volume passing through link at any specific time.

Example Case

There are 4 nodes a,b,c,d. Nodes b,c,d each route traffic only to node a. A unit amount of load is sent from b,d to a. e amount of load sent from c to a.

Link cost = amount of load carried on link

Initial routing :-



(Routes from Src to dest which are in use currently marked by blue.)

We can observe at this instant,

Costs of

- clockwise path from C to a = 1

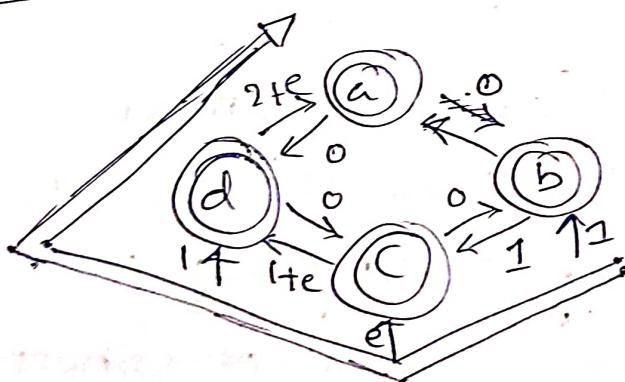
- Counter clockwise path from C to a = $1 + 2e$

\Rightarrow C's least cost path to a is clockwise

There occurs a need to change our initial route.

In a similar way, b's least cost to a is also clockwise (1) rather than currently selected Counter clockwise ($1 + e$)

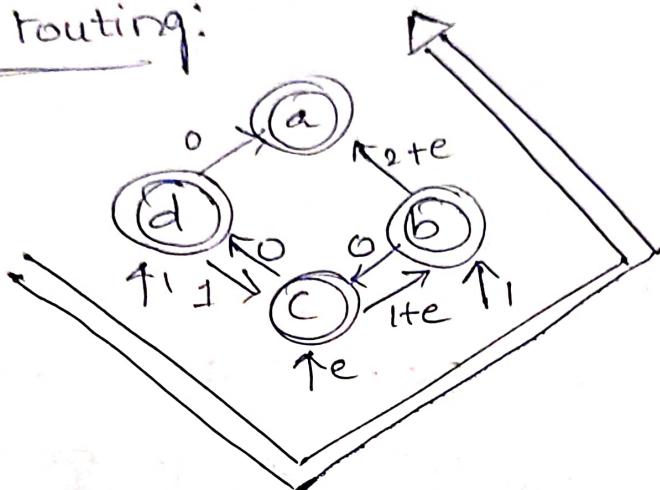
Modified routing



In the next run of algo, b will detect zero cost path to a in anticlockwise direction & corresponding the traffic is routed along minimum cost direction

Now, on the next run

Modified routing:



Here nodes b, c, d detect 0 Cost path to a in clockwise direction and traffic is again routed. Same as before fig. These modifications Oscillate from clockwise to anti-clockwise / vice-versa as long as routing need to be done.

This exact Scenario is known as Oscillation problem in traditional Computer network due to this many routing algorithms do not define link costs based on short term load levels. Software defined networking has been a major advancement which can be used to Overcome Oscillation Problem.

(10) A:

Motivation to adopt Software defined network

Over traditional network

The reason is that the traditional networking is rooted in fixed-function network devices, for example a Switch or route. If the networks functions are implemented as hardware constructs then its speed is usually bolstered.

The recurring hurdle for traditional networks is flexibility. Few APIs are exposed for provisioning software, but this software can't be quickly modified as needed. Traditional networking consist of traits.

- (i) Functions of traditional networking are primarily implemented from dedicated devices using One or more Switches, as well as routers and application delivery controllers.
- (ii) Functionality of traditional networking is largely implemented in dedicated hardware, such as application-specific

Integrated Circuits. One of the negative aspects of this traditional hardware-centric networking is its limitation.

SDN has advantage of generating the framework that bolsters data-intensive applications such as Big data & Virtualization.

In addition to these benefits here are some other reasons data enterprises network management, SDN offers following succinct advantages:

- * Traffic programmability
- * Greater agility
- * Capacity to generate policy-driven network supervision
- * Ability to implement network automation.

Here are the primary advantages that SDN offer:

- * Centralized network provisioning.
- * Holistic enterprise management
- * more granular Security
- * Lower Operating Cost and hardware savings & reduced Capital expenditure.
- * Cloud abstraction
- * Consistent & timely Content delivery.

In addition to these benefits here are some other reasons why data enterprises are choosing SDN

- * SDN lets network operators and administrators adjust their bandwidth as needed
- * Traditional networks are rigid, making it difficult for network operators & administrators to customize the programming of their networks.

Traditional network architectures that are too old, rigid & expensive to scale are out of alignment with today's hybrid cloud combination of traditional, public and private cloud and IT as a Service (ITaaS) deployments. Gives great innovations & reductions in complexity.

For traditional network infrastructure, each switch determine where traffic goes and then directs the traffic based of these determinations. with SDN infrastructure, the process of determination and direction has been decoupled. Switches still direct the traffic how process of determining where the traffic goes is performed automated programmable interface called controller.

This platform allows network administrators to change network-wide settings with a Centralized Console. Whereas traditional network infrastructure might warrant deploying network changes in a piecemeal fashion for individual devices.

SDN infrastructure gives network administrators the flexibility to change network traffic and enables network resources deployment that scales at the same speed as server and storage, redirecting it as needed.

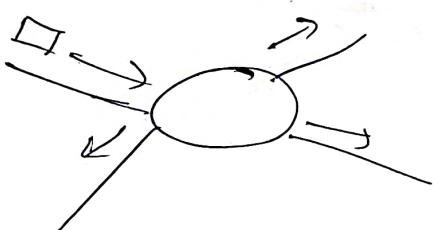
Additionally, the SDN Controller reduces complexity and enables the network to scale as needed.

→ The concept of networks substrate comes when we talk about SDN. It is used for being able to mimic the success in Software industry.

Characteristics include

- (i) Has Simple Common Substrate
- (ii) Build OS on top of hardware which enables easy deployment of networking applications.

Router Example



Basic job of router is receiving packets, checking routing table, forwarding packets out. For building routing table, router has to understand BGP, OSPF, RIP etc.. we can use network substrates to get the routing table from somewhere else (centralized)

This can also be seen in Separate data and control plane based architecture of SDN. Since Control plane is centralized and its on top of all software, hardware it is a network substrate.

(11) A:

Software - Defined Networking

SDN is the physical Separation of the network Control plane from the forwarding plane and where a Control plane Controls Several devices. SDN makes networks agile & flexible. It provides better network Control for cloud Computing Service Providers to respond quickly to ever changing business requirement & preferred over traditional architecture.

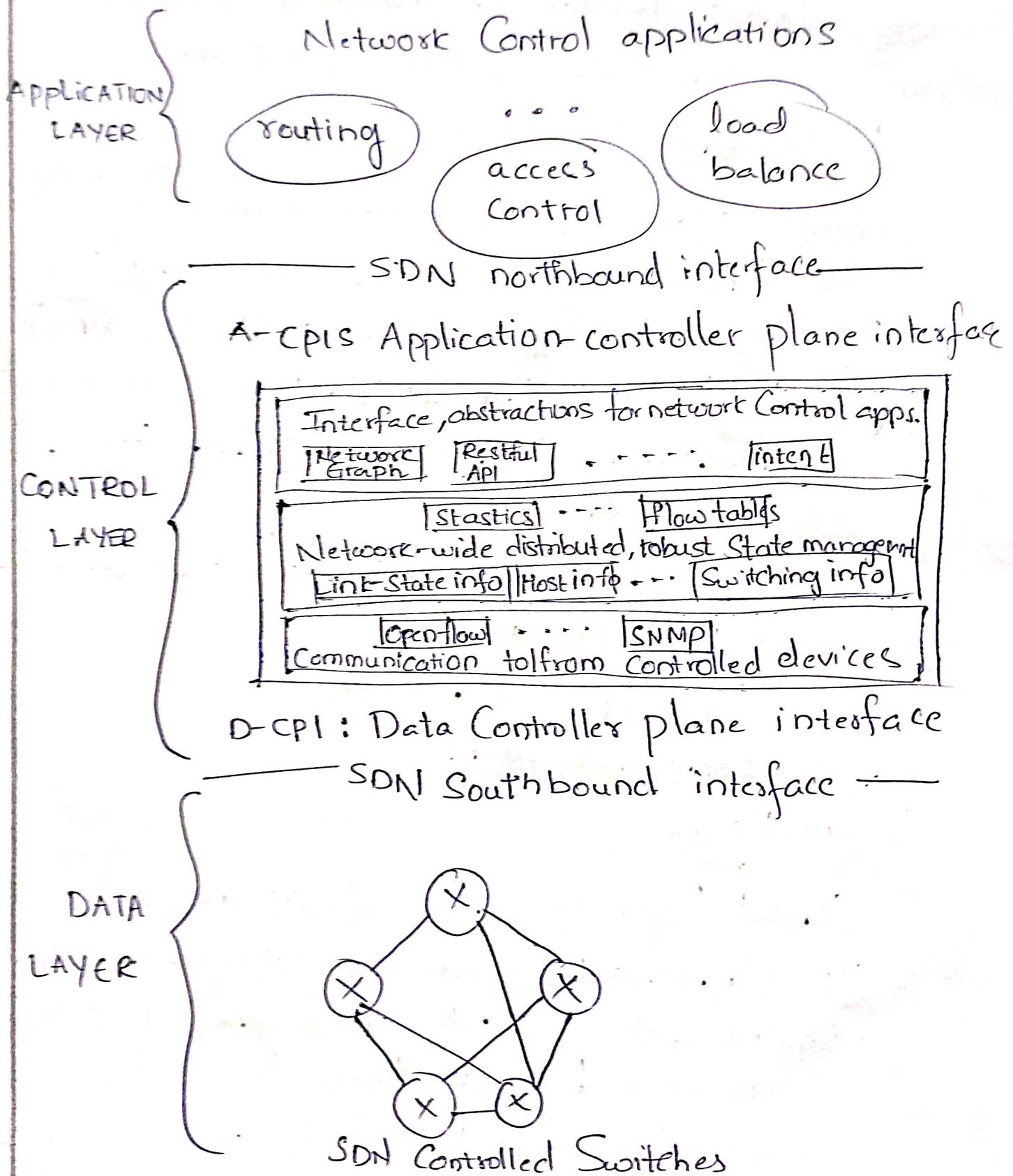
SDN has 3 layers

application layer, Control layer, Infrastructure layer.

Data plane. Consists of network elements, which expose their Capabilities to the Control plane via South bound interface

Control plane sits in middle & has following functions

SDN Architecture diagram



- translate application's requirements and exerts low-level control over network elements
- provide network information to applicants
 - Orchestrate different applications.

Data plane

Data plane consists of data sources and sinks. It functions as a traffic forwarding/processing engine and may have the ability to handle protocols like ARP, LLDP. It has interfaces communicating to the control plane. Core functionality of data plane is

- programmatic control of all functions offered by network element.
- capability advertisement
- event notification.

Control plane

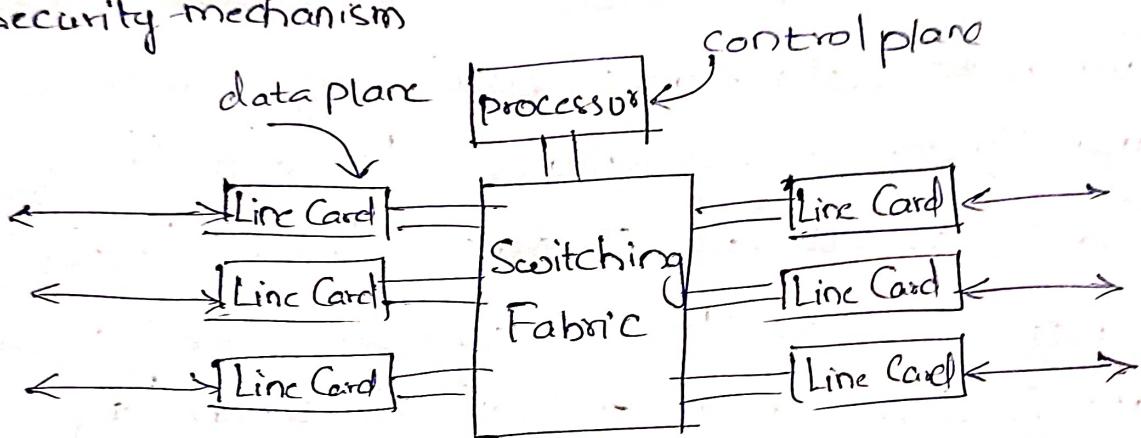
Control Layer is considered as the brain of SDN. It communicates with the data plane using Southbound APIs. The intelligence to this layer is provided by Centralized SDN Controller Software. This controller resides on a server &

manages policies and flow of traffic throughout network.

Core functionality of Control plane is

- Topology and network State information
- Device discovery
- Path Computation.

Security mechanism



Interaction b/w Data & Control plane

Application Plane

SDN Application layer Contains network applications or functions Such as intrusion detection System , load balancing or firewalls . It Communicates with Control layer using northbound API .

Applications Specify The resources and behaviors required from Network, with Context of business & policy agreement. Applications are developed by Specific programming languages & may need to orchestrate multiple - Controllers to achieve the objectives.

Data Center transport impairments & requirements

(12) AE

The shallow packet buffers cause 3

Specific performance impairments

1. Incast

If many flows converge on same interface of a switch over short period of time, packets may exhaust either switch memory or maximum permitted buffer resulting in data losses. This traffic pattern arises naturally from use of partition / aggregate pattern as request for data creates incast at queue of switch port connected to aggregator.

We capture incast instances via packet-level monitoring. Since size of each individual response is small, each individual response leads to loss of packet almost invariably results in a TCP timeout. Thus, whenever timeout occurs response almost always misses aggregator's deadline.

2. Queue buildup

Long-lived, Greedy TCP flows will cause length of bottleneck queue to grow until packets are dropped resulting in familiar saw tooth patterns.

Two impairment occurs when short & long flows traverse same queue. Even when no packets are lost, short flows experience increased latency as they are in queue behind packets from large flows, thus queue buildup impairment occurs & this traffic pattern occurs more frequently.

3. Buffer pressure

It is very common for short flows on one port to be impacted by activity on any of many other ports. Indeed, the loss rate of short flows in this traffic pattern depends on no. of long flows traversing other ports.

The long, Greed TCP flows build up

- queues on their interfaces. Since buffer space is a shared resource, the queue

build up reduces the amount of buffer space available to absorb bursts of traffic from partition / Aggregate traffic. We term this impairment buffer pressure. The result is packet loss and timeouts as in incast but without requiring synchronized flows.

Tension between requirements

High Throughput

High Burst tolerance

Low latency

Deep buffers

Reduced RTT_{min}

> Doesn't help latency

Shallow buffers

AQM - RED

> Avg queue not fast enough for incast

Case Study

Microsoft Bing

- > Measurements from 6000 Server production cluster
- > Instrumentation Passively Collects logs
 - Application-level
 - Socket-level
 - Selected-packet-level

- > More than 150 TB of compressed data
Over a month

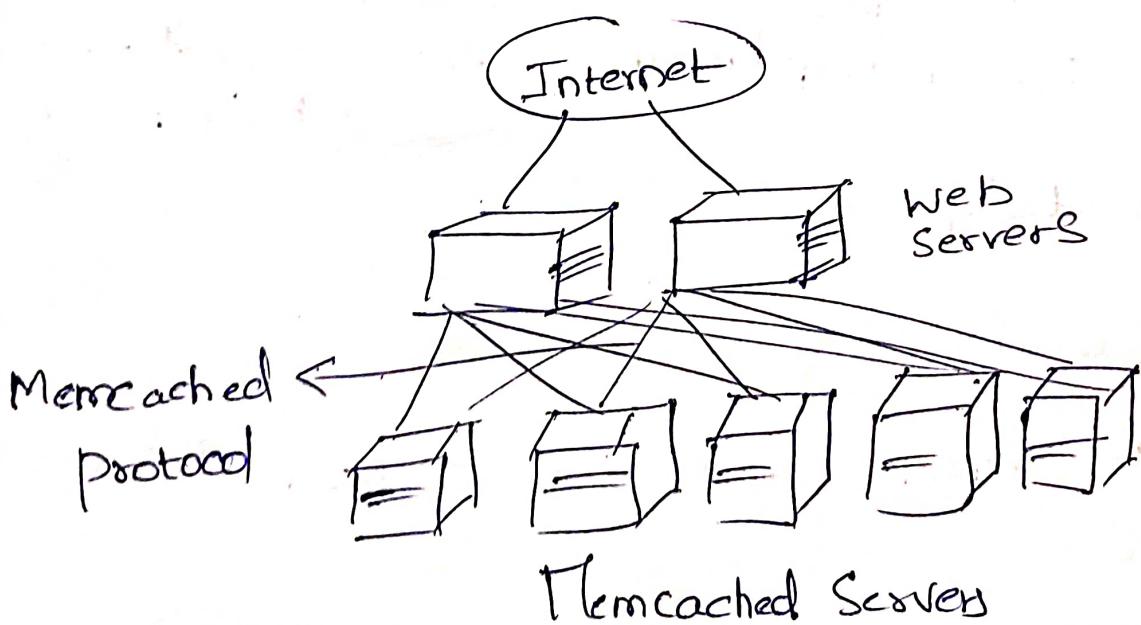
Where as in the case of facebook, it uses

- Partition / Aggregate

Partition / Aggregate - Foundation for many large scale web apps

Facebook

- > Uses partition / Aggregate ~ Multilevel
Aggregators: webservers
Workers: Memcached Servers



Workloads are (facebook)

- > Partition / Aggregate → Delay Sensitive
(query)
- > Short messages → Delay Sensitive
[50 KB, 1 MB]
(coordination, control state)
- > Large flows [1 MB - 50 MB] → Throughput
(Data update) Sensitive

(13) A:

a) Forwarding Abstraction:

purpose: Abstract away forwarding hardware while still be able to express how and where to forward a packet.

The first abstraction relates to how network elements forward packets. At a high enough level of abstraction all network elements perform the same task.

Abstraction I: packet forwarding as a Computational problem

The function of any network element is to

- (i) Receive a packet
- (ii) Observe packet fields
- (iii) Apply algorithms
- (iv) Optionally edit the packet
- (v) Forward or discard the packet

Features:

⇒ Flexible

- forwarding behaviour Specified

by Control plane

- Built from basic set of forwarding primitives

→ Minimal

- Streamlined for Speed and low power

- open and not vendor specific

Openflow is an example of forwarding abstraction

Openflow forwarding abstraction

- Protocol independent

(i) Construct ethernet, IPv4, VLAN, MPLS

(ii) Construct new forwarding methods

Backward Compatible

(i) Run in existing networks

- Technology independent

(i) Switches, Routers, WiFi APs

(ii) Cellular base stations

(iii) WDM / TDM Circuits.

b) Network functions Virtualization

It is a way to virtualize network services, such as routers, firewalls, and load balancers, that have traditionally been run on proprietary hardware. These services are packaged as virtual machines on commodity hardware, which allows service providers to run their network on standard servers instead of proprietary ones. It is one of the primary components of a telco cloud, which is reshaping the telecommunication industry.

With NFV, we don't need to have dedicated hardware for each network function. NVF improves scalability and agility by allowing service providers to deliver new network services and applications on demands without requiring additional hardware resources.

NVF architecture:-

- * Virtualized network functions: Those are software applications that deliver network

functions such as file sharing, directly services, and IP configuration.

* Network function Virtualization infrastructure: It consists of the infrastructure components - compute, storage, networking.

* Management, automation and network orchestration:

It provides a framework for managing NFV and VNF's.

Benefits of using NFV:

With NFV, service providers can run network functions on standard hardware, instead of dedicated hardware.

NFV gives providers the flexibility to run VNFs across different servers or move them around as needed when demand changes. This flexibility lets service providers deliver services and apps faster. For example, if a customer requests a new network function, they can spin up a new VM to handle that request. If the function is no longer needed, the VM can be decommissioned. This can also be a low-risk way to test the value of a potential new service.

c) Service Function Chaining

It is a capability that uses SDN capabilities to create a chain of connected network services, such as L4-7 services like fire walls, network address translation (NAT) and intrusion protection. Network operators can use this to set up suites or catalogs of pathways for traffic to travel through. Any one path can consist of any combination of connected services, depending on the traffic requirements. Different traffic requirements might be more security, lower latency, or an overall high quality of service (QoS). The primary advantage of network service chaining is to automate the way virtual network connections can be setup to handle traffic flows for connected services. For example, an SDN controller could take a chain of services and apply them to different traffic flows depending on the source, destination, or type of traffic.

The chaining capability automates what

traditional network administrators do when they connect up a series of physical L4-7 devices to process incoming and outgoing network traffic, which traditionally may require a lot of manual steps. Service chaining is being included in many SDN and network function virtualisation use cases and deployments, including data中心. Another benefit is when used in combination with SDN, is optimizing the use of network resources and improving application performance. SDN analytics and performance tools can use the best available network resources and help negotiate around network congestion issues.

Service chaining is facing a new networking technology that is purposefully trying to get away from service chaining. That technology is called SASE. The chain in service chaining represents the service that can be connected across the network using software provisioning. This is especially important in NFV world.

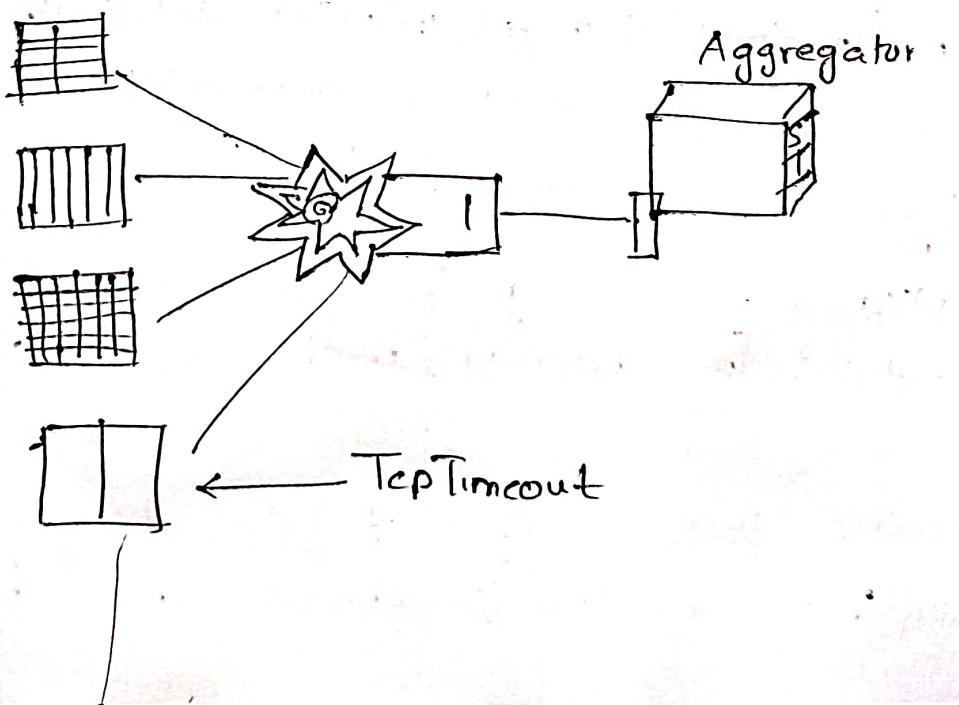
Navya Sree
1801CS06
CS358
NLF

The Technology's capabilities mean that a large number of virtual network functions can be connected together in a NFV environment. Because it's done in softwares using virtual circuits, these connections can be set up and torn down as needed with service chain provisioning through the NFV.

(14) A:

Incast problem in data Centres

TCP-Incast problem refers to TCP performance degradation when it runs in many-to-one traffic pattern. Initially Catastrophic collapse is observed when more no. of Senders Concurrently transmits data blocks to a Single receiver and any Sender is not allowed to send the next block until all others finish transmitting. Therefore in the delay-sensitive online application which employ partition/aggregation Computing. Very few Senders extended the time, so the System responses are included from final result due to missing aggregators deadline.



Solutions recommended to solve TCP incast problem

1) Adjusting System parameters

It improves TCP incast problem by adjusting (decrease)
parameters, such as changing block size, enlarging buffer size and increasing capacity.

Limiting the no. of concurrent flows is also away to improve situation.

2) Designing Enhanced mechanisms

- Reduce RTO min to reduce throughput collapse
- Shrinking MTU (max. transmission limit)

3) Replacing loss-based TCP version

- Adjusting Congestion windows properly according to delay information generated by RTT measurement. Slight adjustment of window size is beneficial to protect the buffer

4) New Transmission protocols.

DCTCP Algorithm: It works because

- (i) High Burst tolerance

- Large buffer headroom → burst fit

- Aggressive marking → sources react before packets are dropped.

NavyaSree
1801CS06
CS358
11

(ii) Low latency

Small buffer Occupancy \rightarrow low quality delay

(iii) High throughput

ECN agency \rightarrow Smooth rate adjustments, low variance

(15) A:

Markov Decision problem

These processes are discrete-time stochastic control processes. They provide a mathematical framework for modelling decision making in situations where outcomes are partly random and partly under the control of a decision maker.

Problem: MDP's are modelled as follows

4-tuple (S, A, P_a, R_a) where :

- (i) S is a set of states called state space.
- (ii) A is a set of actions called action space
- (iii) $P_a(s, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$
(probability that action a in state s at time t will lead to state s' at time $t+1$)
- (iv) $R_a(s, s')$ is the immediate reward received after transitioning from state s to s' due to action a .

A policy function π is a probabilistic mapping from State Space to action Space.

Objective:

The Objective of an MDP process is simply to find the Optimal policy ($\pi(s)$) mapping from State Space (S) to action (A) Space.

To do this, we use a function called Value function

$$V(S) = \sum_{S'} P_{\pi(S)}(S, S') (R_{\pi(S)}(S, S') + \gamma V(S'))$$

This value function would return the expected returns of all rewards possible from state S.

And naturally we would want a policy function which would maximize the expected return.

$$\pi_*(S) = \arg \max_a \left\{ \sum_{S'} p(S'|S, a) (R(S'|S, a) + \gamma V(S')) \right\}$$

We can converge to this optimal policy by iterating through 2 - Steps

(i) Value update

(ii) policy update

which are repeated in some order for all states until no further change takes place.

Both recursively update a new estimation of Optimal policy & state value using an older estimation of those values.

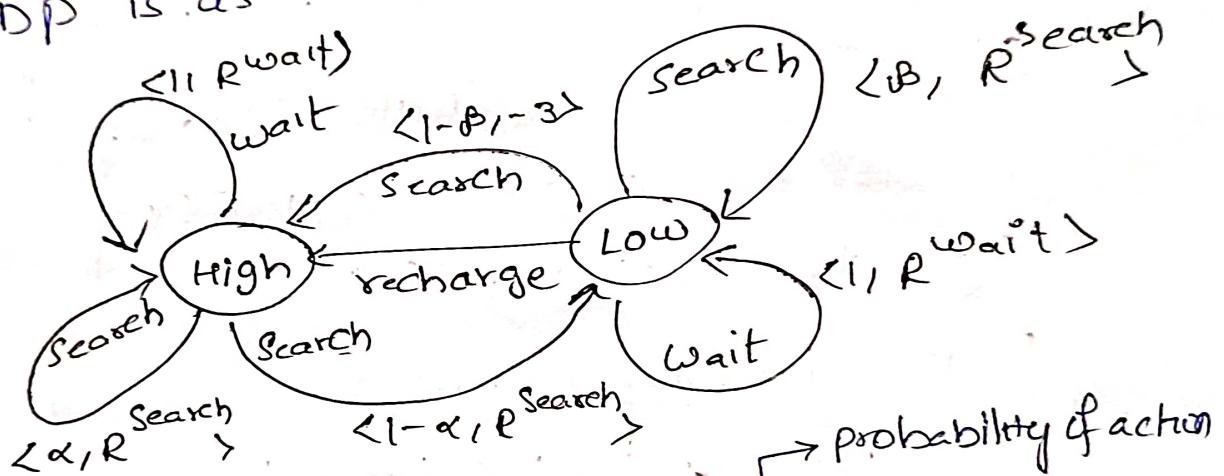
Ex: Consider robot whose purpose is to pick up litter

2 States $S = \{\text{high, low}\}$
Battery Charge

Actions $A = \{\text{wait, Search, Recharge}\}$

- wait ↑
Hold still
- Search ↑
Search for litter
- Recharge ↑
Go to charging station

MDP is as follows,



Edges are represented by $\langle p, r \rangle$ → Reward

Q-Learning

Q-Learning is a model-free RL algorithm to learn value of action in particular state. It does not require a model of environment and it can handle problems with Stochastic transitions & rewards without requiring adaptations.

For any finite MDP (FMDP), Q-learning finds an Optimal policy in the Sense of maximizing expected value of total reward Over all successive steps, Starting from Current state.

* Q-learning Can identify an Optimal action-selection policy for any given FMDP, given infinite exploration time and a partly random policy.

Objective: We estimate Optimal Q-function we assume a Q-table, mapping State action pairs to Values ie., $Q : S \times A \rightarrow R$

Before learning begins, Q is initialized to a possibly arbitrary values. Then at each time 't' the agent selects an action a_t , Observes a reward r_t , enters new state S_{t+1} , Q is updated.

$$Q'(s_t, a_t) = Q(s_t, a_t) + \alpha \left(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right)$$

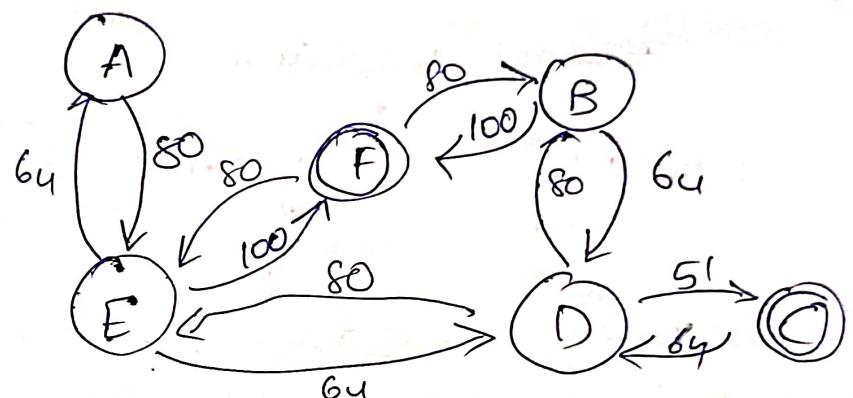
↓ ↓ ↓ ↓
 new value Old value Learning rate discount factor old value
 ↓
 temporal difference

Three factors,

- i) $(1-\alpha)Q(s_t, a_t) \rightarrow$ Current Val weighed by α
- ii) αr_t : Reward if a_t is taken when
- iii) $\alpha \max Q(s_{t+1}, a)$: maximum reward that could be obtained from s_{t+1}

Ex: Consider

following environment



Initial Q table,

$$Q = \begin{bmatrix} & A & B & C & D & E & F \\ A & 0 & 0 & 0 & 0 & 0 & 0 \\ B & 0 & 0 & 0 & 0 & 0 & 0 \\ C & 0 & 0 & 0 & 0 & 0 & 0 \\ D & 0 & 0 & 0 & 0 & 0 & 0 \\ E & 0 & 0 & 0 & 0 & 0 & 0 \\ F & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Final Q table, $Q =$

$$\begin{bmatrix} & A & B & C & D & E & F \\ A & - & - & - & - & 80 & - \\ B & - & - & - & 64 & - & 100 \\ C & - & - & - & 64 & - & - \\ D & - & 80 & 51 & - & 80 & - \\ E & 64 & - & - & 64 & - & 100 \\ F & - & 80 & - & - & 80 & 100 \end{bmatrix}$$

MDP

- (i) Non-temporal difference
- (ii) Discrete Time Stochastic Control process
- (iii) Model based / policy based algorithm
- (iv) utilizes Dynamic Programming
- (v) uses value function & policy function for optimization
- (vi) Used when R, P are Known beforehand

Q-Learning

- (i) Temporal difference
- (ii) Continuous Time Stochastic control processes
- (iii) Model free Algorithm
- (iv) utilizes Bellman Equation
- (v) Utilizes Q-table for Optimization
- (vi) Used when R, P are unpredictable

\Rightarrow MDP performs better than Q-Learning in non temporal environments i.e., when episodes end at a time stamp

\Rightarrow MDP does not suffer from bias, as each update is made using a true sample of what $Q(s,a)$ should be.

(16) a:

Exploration and Exploitation in Q-Learning

The Q-Learning algorithm does not specify what the agent should actually do. The agent learns a Q-function that can be used to determine an optimal action.

There are 2 things that are useful for agent to do

Exploit: The knowledge that it has found for the current states by doing one of the actions 'a' that maximizes $Q[s, a]$

Explore: In order to build a better estimate of optional Q-function.. i.e) It should select a different action from the one that it currently thinks is best

Tradeoff: During learning, it might be the case that the agent chooses the action value which has the highest estimate (greedy action) or something else (non-greedy action).

Doing only one of these doesn't help in Convergence
Key is balancing b/w exploitation & exploration.
Only exploitation might be good in short run but
might miss long term optimum. ∵ even though
initially exploration gives bad results it will
help in convergence

Few algorithms that balance b/w both are

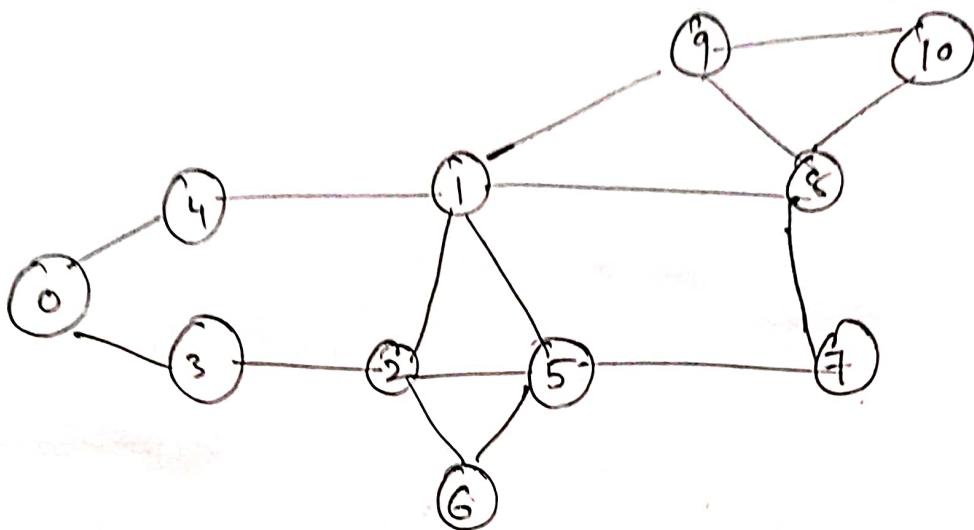
1. Greedy Algorithm
2. Upper Confidence bound
3. Thompson Sampling

Note: The property that needs to be held by
exploitation is "Degradation". Exploitation must
happen in the initial stages of learning &
must die down eventually. or else
it accounts to high Variance

Eg: Decaying E Greedy policy

(17) Ans:

Given Graph,



It is also given, Optimal π -matrix

$$\text{As we know, } \pi^*(s) = \arg \max_a (\pi(s, a))$$

$$\Rightarrow \pi^*(0) = 4 \quad (\text{138 is max val in 0 row})$$

$$\pi^*(4) = 1 \quad (174 \text{ " " " } 4 \text{ row})$$

$$\pi^*(1) = 8 \quad (218 \text{ " " " } 1 \text{ row})$$

$$\pi^*(8) = 10 \quad (274 \text{ " " " } 8 \text{ row})$$

\therefore Path $\Rightarrow 0 \rightarrow 4 \rightarrow 1 \rightarrow 8 \rightarrow 10$

also possible path is $0 \rightarrow 4 \rightarrow 1 \rightarrow 9 \rightarrow 10$