# Open Items for Flooding Reduction

Huaimo Chen
Dean Cheng
Mehmet Toy
Yi Yang
Aijun Wang
Xufeng Liu
Yanhe Fan
Lei Liu

# Overview

➢ FT Consistency Check: Behavior for Inconsistency

➢ Flooding Negotiation (FN) bit

➢ Transfer between Flooding Reduction and Normal Flooding

➢ Enhancement related to Area Leader Sub-TLV

➢ Enhancements on FT Encoding

➢ Backup Paths for FT Partitions

# FT Consistency Check: Behavior for Inconsistency (1)

- LEEF bit added into the draft based on the FT bit for a link. The bit advertised by one end node of the link indicates whether the link is on the flooding topology

A

LEEF-bit=1

B

LEEF-bit=0

Should we add Behavior for Inconsistency below?
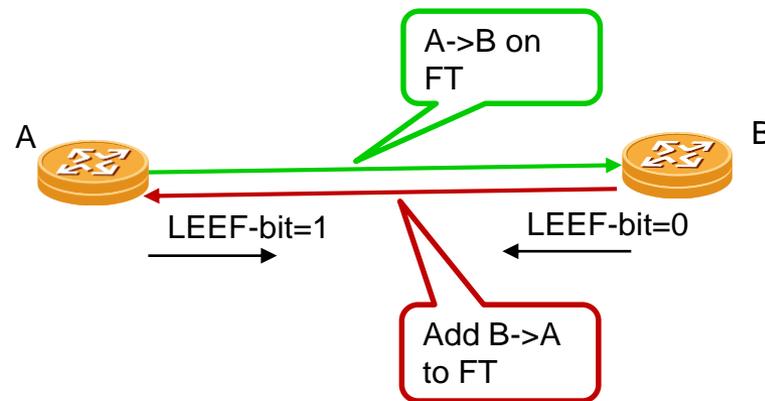
For link L between A and B, if the LEEF bit for link L advertised by A is different from the one by B, then the flooding topology is inconsistent, the node receiving the LEEF bit set to one for link L from the other node adds link L on the FT temporarily.

If one of the two nodes receives the LEEF bit set to one for link L from the other, but advertises the LEEF bit set to zero for link L for a given time such as 5 seconds, then a warning is issued or logged.

# FT Consistency Check: Behavior for Inconsistency (2)

Reasons for adding behavior:
- May fix problem of the FT quickly
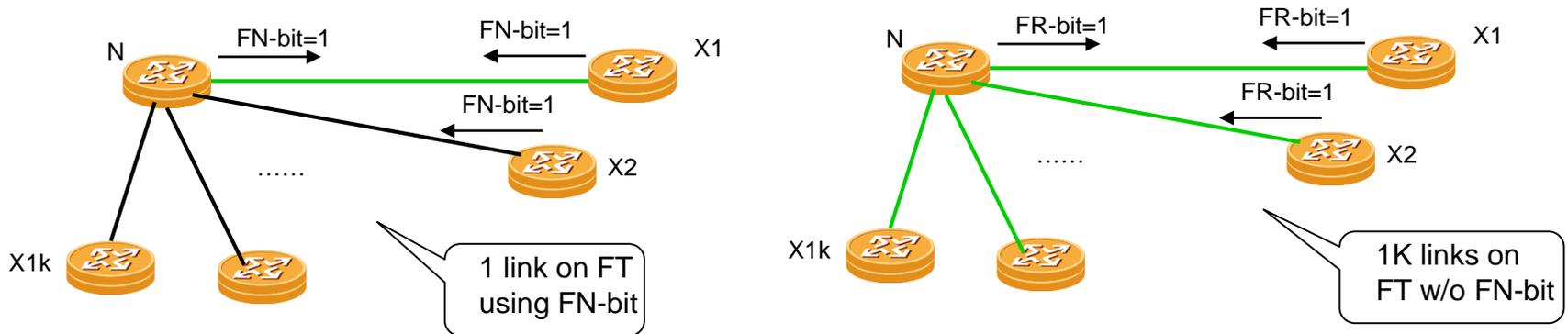- Better to give a Warning by node (such as B) that detects inconsistency

A->B on FT

A        B

LEEF-bit=1        LEEF-bit=0

Add B->A to FT

Reasons for not adding behavior:
- No change to flooding behavior
- Tools to detect inconsistency

# Flooding Negotiation (FN) bit (1)

**Problem to be solved:**
If node N reboots with 1K links, which one or few links should be added to FT temporarily? Adding all is likely to trigger a cascade failure.



Xi adds link Xi-N to the FT temporarily after it sends and receives the FN bit=1 to/from N;
N adds a selected link to the FT temporarily after it receives and sends the FN bit=1 from/to Xi.

During the process of the adjacency establishment between Xi and N, Xi sends a FN-bit=1 in its Hello to N, N selects one link/node (or a few links) for temporarily flooding and sends only to this selected node a FN-bit=1 in its Hello.

There are two different cases in which a link is to be added to the FT temporarily.
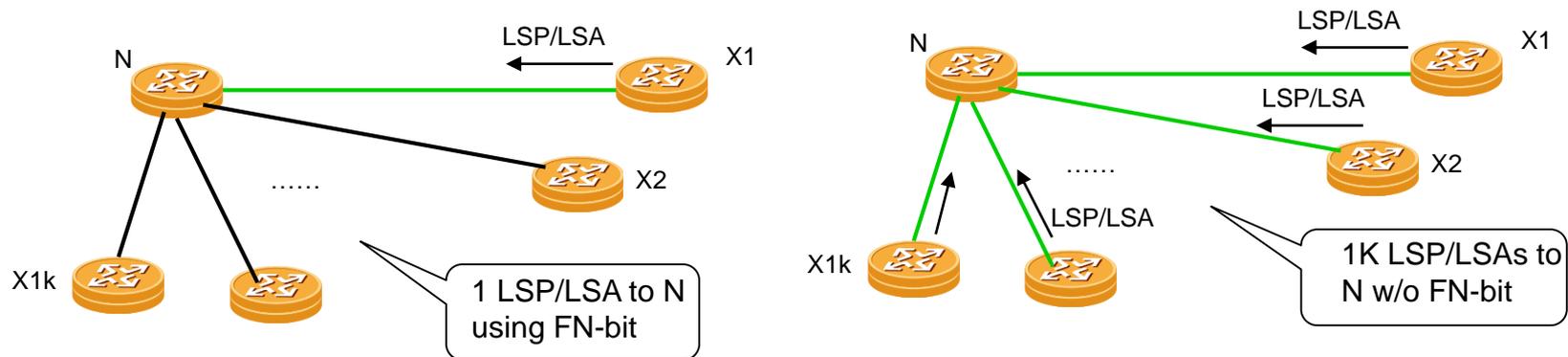In one case, a negotiation is needed to be done before a link is to be added to the FT temporarily.
In the other case, no negotiation is needed. It is determined that a link is added to the FT temporarily.

# Flooding Negotiation (FN) bit (2)

Reasons for adding FN-bit:
- Resolves the problem (Adding all is likely to trigger a cascade failure)
- It is the objective of flooding reduction to add one or two links to the FT.
- If 1K links are added to the FT after 1K adjacencies are up, for one updated LSP/LSA, 1K LSP/LSAs of same copy will be flooded to node N from X1,X2, ..., X1k.
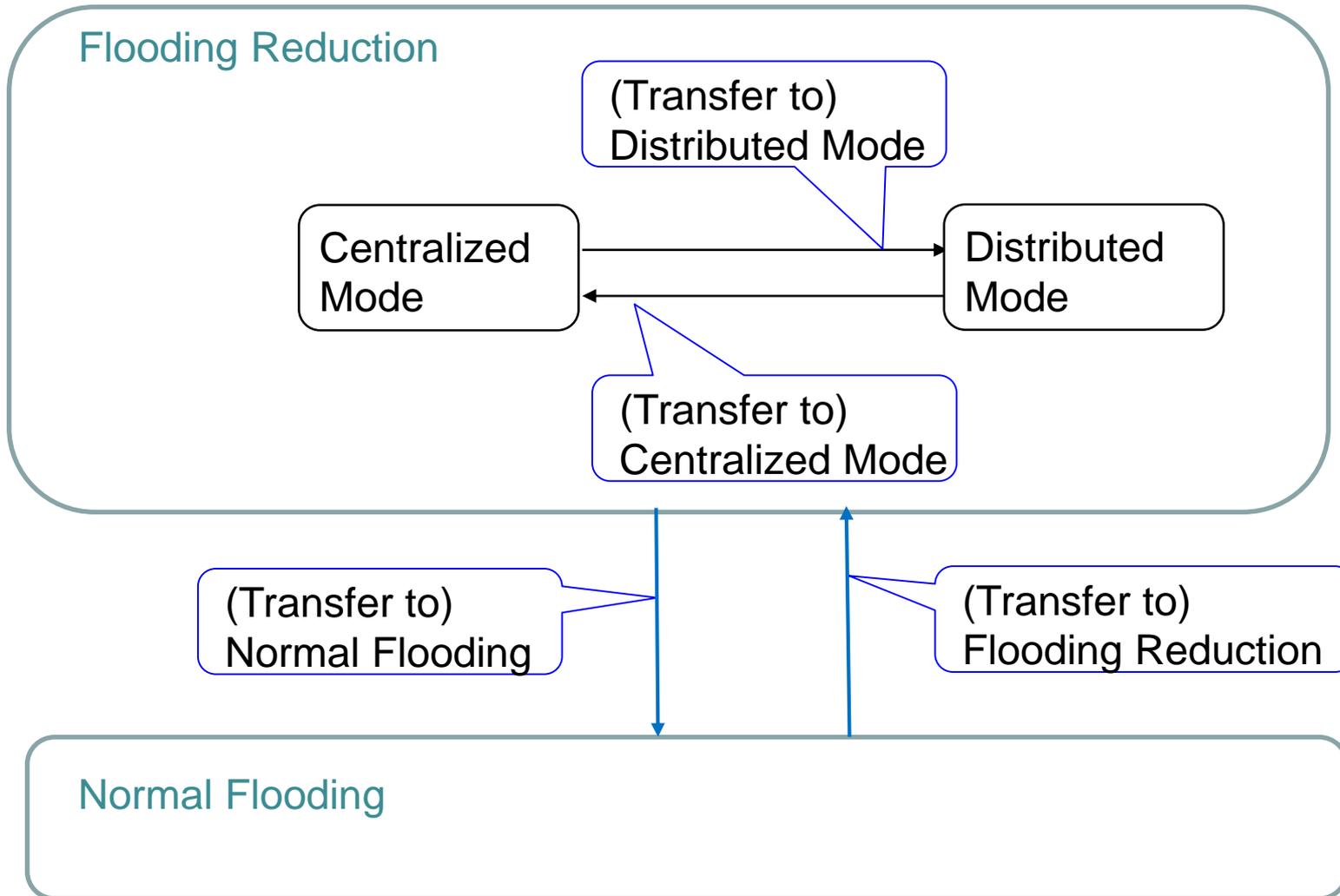


Reasons for not adding FN-bit:
- Problem is resolved by limiting the rate of adjacency establishments

Flooding Reduction <==> Normal Flooding Quickly, Easily and Smoothly

**Flooding Reduction**

(Transfer to)
Distributed Mode

Centralized Mode → Distributed Mode

(Transfer to)
Centralized Mode

(Transfer to)
Normal Flooding

(Transfer to)
Flooding Reduction

**Normal Flooding**

Flooding Reduction ==> Normal Flooding:
1. On leader, "(Transfer to) Normal Flooding" is configured;
2. Leader advertises "(Transfer to) Normal Flooding" to all nodes;
3. Every node transfers to Normal Flooding after obtaining instruction.

- For centralized mode, after transferring to Normal Flooding, leader stops computing and advertising FT, each of the other nodes stops receiving and building FT.
- For distributed mode, every node stops computing and building FT.

At this point, IGP in area has transferred to Normal Flooding from Flooding Reduction (either centralized mode or distributed mode).

Normal Flooding  ==> Flooding Reduction:

For centralized mode (i.e., when centralized mode is configured),
1.  Leader advertises "Flooding Reduction" in the centralized mode to all the other nodes
2.  Leader computes FT and advertises FT to the other nodes;
3.  Each node floods link states using FT after it receives/has whole FT.

For distributed mode (i.e., when distributed mode is configured),
1.  Leader advertises "Flooding Reduction" in the distributed mode including the algorithm to all the other nodes;
2.  Each node computes its FT and floods the link states using the FT.

At this point, IGP in area has transferred to Flooding Reduction (either centralized mode or distributed mode) from Normal Flooding.

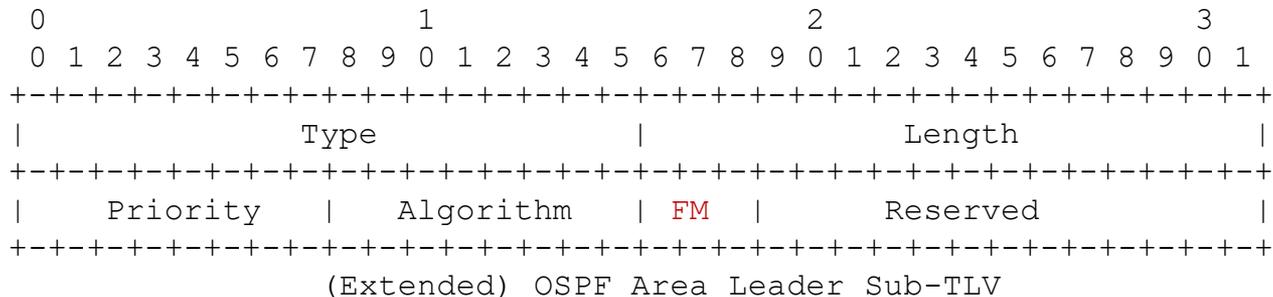# Transfer between Flooding Reduction and Normal Flooding (4)

To support the above behaviors, in one Explicit way, Area Leader Sub-TLV is extended. Three bits of one octet may be used to indicate a flooding mode (FM) such as "Normal Flooding" or "Flooding Reduction". The other bits are reserved. The values proposed for FM are:
- 1 (or 0) for "Flooding Reduction" (centralized or distributed mode is implied/indicated by the algorithm)
- 2 for "Normal Flooding"

```
For OSPF Area Leader Sub-TLV,
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              Type             |              Length           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|    Priority   |   Algorithm   |            Reserved           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
               (Current) OSPF Area Leader Sub-TLV
is extended to
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              Type             |              Length           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|    Priority   |   Algorithm   | FM  |          Reserved       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
               (Extended) OSPF Area Leader Sub-TLV
FM = 1 (or 0):  Flooding Reduction
        Alogrithm  =  0:  Centralized Mode; Algorithm = N (N>0): Distributed
Mode.
FM = 2: Normal Flooding
```

To support the above behaviors, in one Implied way, in centralized mode

1.  "(Transfer to) Normal Flooding" is configured on the leader,
2.  Leader advertises its Area Leader Sub-TLV with Algorithm=0 and stops computing FT and advertising FT, which implies to all nodes "Normal Flooding" .
3.  Every node is still doing Flooding Reduction until FT ages out
4.  After 1 hour for OSPF, FT ages out, every node transfers to Normal Flooding.

To speed up transfer, at step 2, leader needs flush FT.

For transferring Flooding Reduction in distributed mode to Normal Flooding in Implied way, it seems complicated.

At step 3 above, every node will do Flooding Reduction in distributed mode forever even after it receives Area Leader Sub-TLV with Algorithm=0. The Sub-TLV is the same as the one indicating/instructing transfer from distributed to centralized mode. Every node continues doing flooding reduction in distributed mode until it receives FT.

To support the above behaviors, in one Implied way, in distributed mode

1.  "(Transfer to) Normal Flooding" is configured on the leader,
2.  Leader advertises its Area Leader Sub-TLV with Algorithm=0 and starts computing FT and advertising FT, which transfers to centralized mode.
3.  Every node transfers from distributed mode to centralized mode after receiving FT.
4.  Leader stops computing FT and advertising FT, which implies to all nodes "Normal Flooding".
5.  Leader flushes FT.
6.  Every node transfers to Normal Flooding after FT ages out.

For transferring Flooding Reduction in distributed mode to Normal Flooding in Implied way,
Every node needs to transfer from distributed mode to centralized mode, does flooding reduction in centralized mode for a short time, and then transfers to Normal Flooding from Flooding Reduction.

Reasons for Flooding Reduction <==> Normal Flooding (in Explicit way):
• Transfer between them Quickly, Easily and Smoothly

Reasons for not Flooding Reduction <==> Normal Flooding (in Explicit way):
• Flooding Reduction ==> Normal Flooding in Implied way

Flooding Reduction ==> Normal Flooding in Implied way
• Transferring Flooding Reduction in distributed mode to Normal Flooding in Implied way is complicated. Every node needs to transfer to centralized mode from distributed mode and then to Normal Flooding from centralized mode

# Enhancement related to Area Leader Sub-TLV (1)

The leader advertises an Area Leader Sub-TLV to indicate/instruct whether centralized or distributed mode is to be used by all the nodes in the area
(Algorithm = 0: centralized mode; Algorithm = N (N>0): distributed mode and N is the algorithm to be used by every node to compute flooding topology).
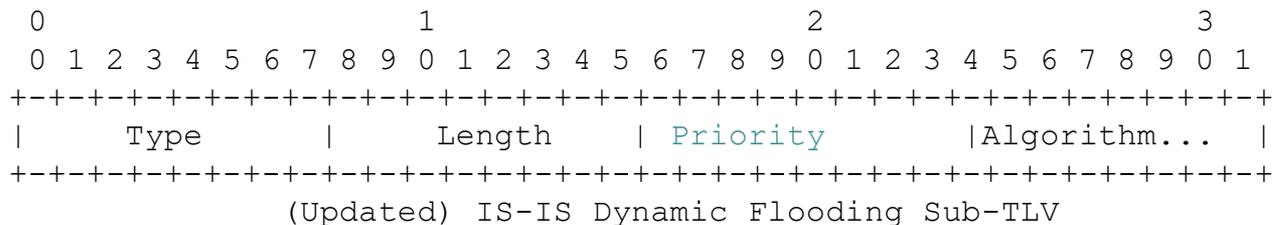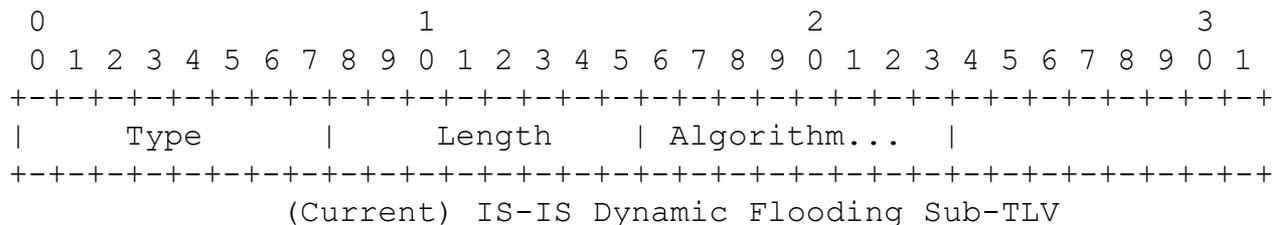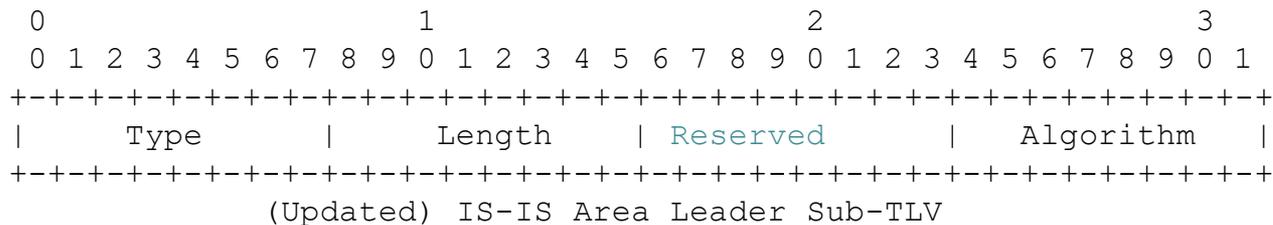
```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Type      |     Length    |   Priority    |   Algorithm   |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
           (Current) IS-IS Area Leader Sub-TLV
```

- Each of multiple nodes (up to all nodes) advertises an Area Leader Sub-TLV to indicate/instruct whether centralized or distributed mode is to be used and to announce its priority for becoming the leader. Hard to understand
- Only leader needs to advertise an Area Leader Sub-TLV to indicate/instruct which mode is to be used by all nodes.

Proposes enhancement:
- Move priority from Area Leader Sub-TLV to Dynamic Flooding Sub-TLV
- Only leader advertises an Area Leader Sub-TLV

5-1

# Enhancement related to Area Leader Sub-TLV (2)

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Type      |     Length    |   Reserved    |   Algorithm   |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
          (Updated) IS-IS Area Leader Sub-TLV


 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Type      |     Length    | Algorithm...  |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
          (Current) IS-IS Dynamic Flooding Sub-TLV

 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Type      |     Length    |   Priority    |Algorithm...   |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
          (Updated) IS-IS Dynamic Flooding Sub-TLV
```

Similarly for OSPF below

# Enhancement related to Area Leader Sub-TLV (3)

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Type              |             Length            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Priority    |   Algorithm   |             Reserved          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
            (Current) OSPF Area Leader Sub-TLV
is changed to
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Type              |             Length            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Reserved    |   Algorithm   |             Reserved          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
            (Updated) OSPF Area Leader Sub-TLV
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Type              |             Length            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Algorithm ... |                                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
            (Current) OSPF Dynamic Flooding Sub-TLV
is changed to
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Type              |             Length            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Priority    | Algorithm ... |                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
            (Updated) OSPF Dynamic Flooding Sub-TLV
```

5-3

# Enhancement related to Area Leader Sub-TLV (4)

The saving on space after moving should be (N-1)*S1 – N. S1 is the size of Area Leader Sub-TLV. S1 is 4 in IS-IS; S1 is 8 in OSPF.
For N = 1000 and IS-IS, the saving is (1000-1)*4 - 1000 = 2996 (bytes)
For N = 1000 and OSPF, the saving is (1000-1)*8 - 1000 = 6992 (bytes)

Reasons for enhancement:
- Easy to understand.
  - Area Leader Sub-TLV is for instructions to all nodes.
  - Only leader gives instructions. Before enhancement, it seems that multiple nodes give instructions since each of them advertises an Area Leader Sub-TLV containing instructions.
- Saving on space at most: (N-1)*S1 – N. S1 is 4 in IS-IS; S1 is 8 in OSPF.

Reasons for no enhancement:
- Only each of a few nodes advertises an Area Leader Sub-TLV to indicate/instruct whether centralized or distributed mode is to be used and to announce its priority for becoming the leader. Thus moving priority (one byte) from Area Leader Sub-TLV to Dynamic Flooding Sub-TLV may use more space. At worse case, (N – 1) bytes more may be used.

# Enhancements on FT Encoding (1)

For the three encodings below,
Encoding A: Plain Path TLVs, where each node index uses 2 bytes,
Encoding B: Compact Path TLVs, where each of node indexes in Path TLVs has the same size given by a field in the Area Node IDs TLV with L set to one, which is the size (in bits) of the maximum node index value,
Encoding C:  Block TLVs,
it seems that C is more efficient than A and B in general, B is more efficient than A.
For example, for representing 63 nodes flooding topology (FT for short) of a binary tree in IS-IS, some comparisons among three encodings are listed in the table below

| Encoding | Bytes used | Comparisons |
|---|---|---|
| A (Plain Path TLVs) | 248 | 179/248 = 0.72 (72%) |
| B (Compact Path TLVs) | 179 | 121/179 = 0.68 (68%) |
| C (Block TLVs) | 121 | 121/248 = 0.49 (49%) |

From the table, we can see that 248, 179 and 121 bytes are used for representing the FT by A, B and C respectively. 121/248 = 0.49 (49%) means that C uses about 49% of the flooding resource that A uses. 121/179 = 0.68 (68%) means that C uses about 68% of the flooding resource that B uses. 179/248=0.72 (72%) means that B uses about 72% of the flooding resource that A uses.

From this example, we can see that C is more efficient than A and B, and B is more efficient than A.

# Enhancements on FT Encoding (2)

Encoding A: Plain Path TLVs
- Break down FT into Paths
- Each Path is represented by a Path TLV, containing indexes of nodes in Path, each index is two bytes.

Encoding B: Compact Path TLVs.
- Break down FT into Paths
- Each Path is represented by a Path TLV, containing indexes of nodes in Path, each index is represented in index size bits, given by a field in Area Node IDs TLV with L=1.

Encoding C:  Block TLVs
- Break down FT into Blocks
- Each Block is represented by a Block TLV, containing indexes of nodes in Block, each index is represented in index size bits, given by the first field in TLV.

# Enhancements on FT Encoding (3)

A Block:

A local node, the number of adjacent nodes, and the adjacent/remote nodes of the local node. (This part is similar to the one in a router LSA to represent the part of the topology from the local node to the adjacent nodes of the local node, which can be considered as a block of the topology in one level. This block can be extended to multiple levels. In LSA, link is uni-directional. In Block, link is bi-directional.)

Each of the adjacent nodes has an extension flag bit E. An adjacent/remote node with E = 1 is considered as a new local node, and its adjacent nodes are added.

—— Link on flooding topology

Router LSA from LN1:
| | |
|---|---|
| Local Node | LN1 |
| #Adjacent Nodes | 3 |
| Adjacent Nodes | RN1, RN2, RN3 |

| | |
|---|---|
| Local Node | LN1 |
| #Adjacent Nodes | 3 |
| Adjacent Nodes | RN1, RN2, RN3 |

7 nodes Block

| | | | | | |
|---|---|---|---|---|---|
| (as Local Node E=1 | RN1) | | (as Local Node E=1 | RN3) |
| #Adjacent Nodes | 1 | | #Adjacent Nodes | 2 |
| Adjacent Nodes | RN11 | | Adjacent Nodes | RN31,R32 |

# Enhancements on FT Encoding (4)

Reasons for using Block TLVs:
- More efficient than Plain and Compact Path TLVs in general.
- Easy to break down FT into Blocks

Reasons for using Compact Path TLVs:
- More efficient than Plain Path TLVs.

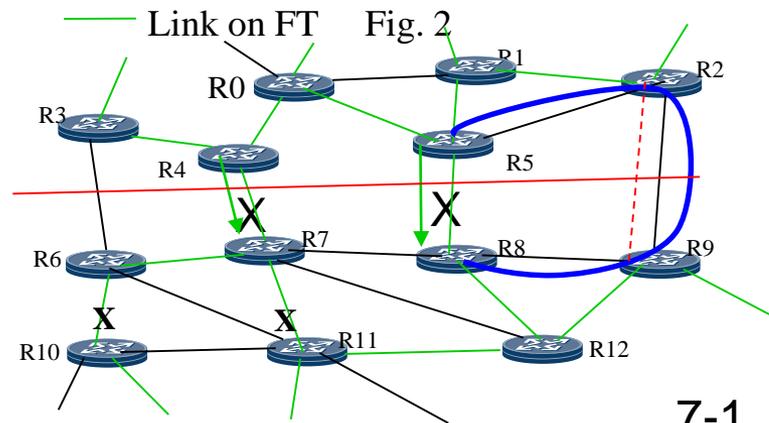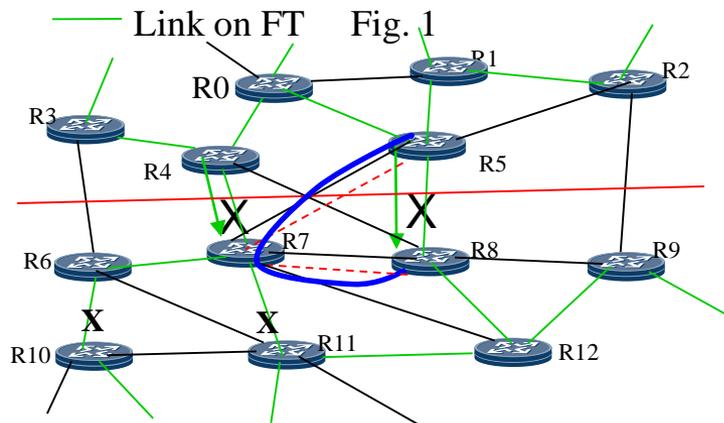Reasons for using Plain Path TLVs:
- Easier to encode Plain Path TLVs.

# Backup Paths for FT Partitions (1)

When multiple link failures on the FT happen, for each failure such as link between A and B down, a unique backup path for the link will be obtained and connect A and B even though the FT is partitioned, A and B are in different partitions, and nodes in different partitions have different LSDBs.

Assume that node A's ID is smaller than node B's ID.
- Node A will compute a minimum hop path from A to B using its LSDB. Suppose that the path is A-X-…-B
- Node X will compute a minimum hop path from A to B using its LSDB.
- The path is the same (i.e., A-X-…-B) if X and A are in the same partition; otherwise (i.e., A and X are in different partitions, link A-X is not on the FT), A will add link A-X to the FT temporarily and the LSDBs of A and X are resynchronized over link A-X; then the backup path computed by A and X will be the same. In this way, the same backup path will be computed by the nodes along the path and enabled for temporary flooding. This path fixes the FT partitions caused by link A-B down.



7-1

# Backup Paths for FT Partitions (2)

For the failure of a node Y on the FT, suppose that there are multiple (i.e., two or more) nodes that were connected to the failed node Y through the links on the FT.

For each pair of these multiple nodes, a unique backup path between this pair is computed and enabled for temporary flooding. Thus the backup paths will connect these multiple nodes on the FT, and the FT partition caused by multiple failures including the failure of node Y is fixed through the backup paths for the failed node Y and the backup paths for the other failures.

# Backup Paths for FT Partitions (3)

Reasons for using Backup Paths:
- A unique backup path between two nodes will be computed, and connect these two nodes on the FT
- FT partitions will be fixed by using almost minimum number of links
- Convergence should be faster since it does not depend on FT computed by leader
- Algorithm for computing a backup path is simple. It is a simplified SPF to find a unique minimum hop path from node A to B.

Reasons for not using Backup Paths:
- Using Rate Limiting Method

Rate Limiting Method
- Nodes which decide to enable temporary flooding also have to decide whether to do so on a subset of the edges which are currently not part of the flooding topology or on all the edges which are currently not part of the flooding topology.  Doing the former risks a longer convergence time as it is possible that the initial set of edges enabled does not fully repair the flooding topology.  Doing the latter risks introducing a flooding storm which destablizes the network.
- It is recommended that a node implement rate limiting on the number of edges on which it chooses to enable temporary flooding.  Initial values for the number of edges to enable and the rate at which additional edges may subsequently be enabled is left as an implementation decision.

# Backup Paths for FT Partitions (4)

Reasons for using Backup Paths:
- It is less likely to have a flooding storm since minimum links are used to fix FT partition.
- Convergence should be faster since it is not through iterations of leader's FTs.
- Algorithm is simple.

Reasons for using Rate Limiting Method:
- Method is simple

# Thanks

# Distributed ==>Centralized Mode

In distributed mode,
1. Configuring centralized mode (or changing the distributed mode to centralized mode through configuration);
2. Leader advertises "Flooding Reduction" in the centralized mode to all the other nodes by set Algorithm = 0 in Area Leader Sub-TLV;
3. Leader computes the flooding topology and advertises the flooding topology to the other nodes;
4. Each node floods the link states using the flooding topology after it receives/has the whole flooding topology.
5. Each node uses the distributed flooding reduction (i.e., floods the link states over its local links on the flooding topology computed and built by itself) until the centralized flooding reduction is fully functional for a given time such as 5 seconds.

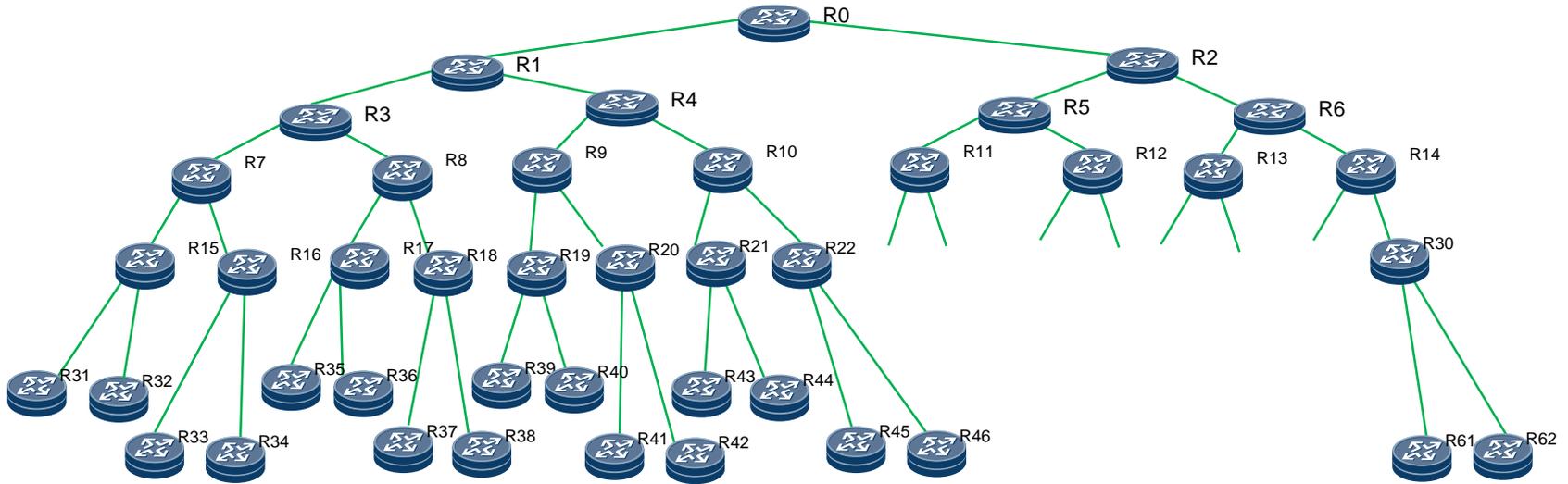# Centralized ==> Distributed Mode

In centralized mode,
1.  Configuring distributed mode (or changing the centralized mode to distributed mode through configuration);
2.  Leader advertises "Flooding Reduction" in the distributed mode to all the other nodes by set Algorithm = n (n > 0) in Area Leader Sub-TLV;
3.  Each node computes its flooding topology and floods the link states using the flooding topology.
4.  Each node uses the centralized flooding reduction (i.e., floods the link states over its local links on the flooding topology computed by the leader of the area) until the distributed flooding reduction is fully functional for a given time such as 5 seconds. After this time, the node stops its centralized flooding reduction. The leader stops computing the flooding topology, advertising it to all the other routers, and using this flooding topology to flood the link states. Each of the other nodes stops receiving and building the flooding topology computed by the leader.

# Election of the DIS

On a LAN, <u>one of the routers elects itself the DIS</u>, <u>based on interface priority (the default is 64).</u> If all interface priorities are the same, the router with the highest subnetwork point of attachment (SNPA) is selected. The SNPA is the MAC address on a LAN, and the local data link connection identifier (DLCI) on a Frame Relay network. If the SNPA is a DLCI and is the same at both sides of a link, the router with the higher system ID becomes the DIS. <u>Every IS-IS router <span style="color:red">interface</span> is assigned both a L1 priority and a L2 <span style="color:red">priority</span> in the range from 0 to 127.</u>

The <u>DIS election is <span style="color:red">preemptive</span></u> (unlike OSPF). If a new router boots on the LAN with a higher interface priority, the new router becomes the DIS. It purges the old pseudonode LSP and floods a new set of LSPs.

# 63 nodes FT of binary tree



Break down FT into paths:
1 path, 11 nodes: R31-R15-…R0-R2-…R62; 2 paths, 2 nodes: R15-R32, R30-R61; 2x(2paths: 2 nodes R16-R33, 3 nodes R7-R16-R34); …