

7 Supervised Overlay Networks I

Every application run on multiple machines needs a mechanism that allows the machines to exchange information. An easy way of solving this problem is that every machine knows the domain name or IP address of every other machine. While this may work well for a small number of machines, large-scale distributed applications such as file sharing or grid computing systems need a different, more scalable approach: instead of forming a clique (where everybody knows everybody else), each machine should only be required to know some small subset of other machines. This graph of knowledge can be seen as a logical network interconnecting the machines, which is also known as an *overlay network*. A prerequisite for an overlay network to be useful is that it has good topological properties. Among the most important are:

- *Degree*: Ideally, the degree should be kept small to avoid a high update cost if a node enters or leaves the system.
- *Diameter*: The diameter should be small to allow the fast exchange of information between any pair of nodes in the network.
- *Node expansion*: The node expansion of a graph $G = (V, E)$ is defined as

$$\beta(G) = \min_{U \subseteq V: |U| \leq |V|/2} \frac{|N(U)|}{|U|}$$

where $N(U)$ is the set of neighbors of U . To ensure a high fault tolerance, the node expansion should be as large as possible.

The question is how to realize such an overlay network in a distributed environment where peers may continuously enter and leave the system. This will be the topic of our investigations for the coming weeks.

We start in this section with the study of *supervised* overlay networks. These networks were investigated, e.g., in [1, 2, 3]. In a supervised overlay network, the topology is under the control of a special machine (or node) called the *supervisor*. All nodes that want to join or leave the network have to declare this to the supervisor, and the supervisor will then take care of integrating them into or removing them from the network. All other operations, however, may be executed without involving the supervisor. In order for a supervised network to be highly scalable, two central requirements have to be fulfilled:

1. The supervisor needs to store at most a polylogarithmic amount of information about the network at any time (i.e., if there are n nodes in the network, storing contact information about $O(\log^2 n)$ of these nodes would be fine, for example), and
2. it takes at most a constant number of communication rounds to include a new node into or exclude an old node from the network.

A *communication round* is over once all the packets that existed at the beginning of the communication round have been delivered. The packets generated by these packets will participate in the next communication round.

We show in the following how these requirements can be achieved, using a general approach called the recursive approach. To simplify the presentation, we assume that all departures are *graceful*, i.e., every node leaving the system informs the supervisor about this and may provide some additional information simplifying the task of the supervisor to repair the network.

7.1 The recursive approach

In the recursive approach, the supervisor assigns a *label* to every node that wants to join the system. The labels are represented as binary strings and are generated in the following order:

0, 1, 01, 11, 001, 011, 101, 111, 0001, 0011, 0101, 0111, 1001, 1011, ...

Basically, when stripping off the least significant bit, then the supervisor is first creating all binary numbers of length 0, then length 1, then length 2, and so on. More formally, consider the mapping $\ell : \mathbb{N}_0 \rightarrow \{0, 1\}^*$ with the property that for every $x \in \mathbb{N}_0$ with binary representation $(x_d \dots x_0)_2$ (where d is minimum possible),

$$\ell(x) = (x_{d-1} \dots x_0 x_d) .$$

Then ℓ generates the sequence of labels displayed above. In the following, it will also be helpful to view labels as real numbers in $[0, 1)$. Let the function $r : \{0, 1\}^* \rightarrow [0, 1)$ be defined so that for every label $\ell = (\ell_1 \ell_2 \dots \ell_d) \in \{0, 1\}^*$,

$$r(\ell) = \sum_{i=1}^d \frac{\ell_i}{2^i} .$$

Then the sequence of labels above translates into

0, 1/2, 1/4, 3/4, 1/8, 3/8, 5/8, 7/8, 1/16, 3/16, 5/16, 7/16, 9/16, ...

Thus, the more labels are used, the more densely the $[0, 1)$ interval will be populated. Furthermore, we will use the function $b : [0, 1) \rightarrow \{0, 1\}^*$ that translates a real number back into a label.

When using the recursive approach, the supervisor aims to maintain the following invariant at every step:

Invariant 7.1 *The set of labels used by the nodes is $\{\ell(0), \ell(1), \dots, \ell(n-1)\}$, where n is the current number of nodes in the system.*

This invariant is preserved when using the following simple strategy:

- Whenever a new node v joins the system and the current number of nodes is n , the supervisor assigns the label $\ell(n)$ to v and increases n by 1.
- Whenever a node w with label ℓ wants to leave the system, the supervisor asks the node with currently highest label $\ell(n-1)$ to change its label to ℓ and reduces n by 1.

How does this strategy help us with maintaining dynamic overlay networks? We will see how this works in the following subsections. To keep things simple, we start with a cycle.

7.2 Recursively maintaining a cycle

We start with some notation. In the following, the label assigned to some node v will be denoted as ℓ_v . Given n nodes with unique labels, we define the *predecessor* $\text{pred}(v)$ of node v as the node w for which $r(\ell_w)$ is closest from below to $r(\ell_v)$, and we define the *successor* $\text{succ}(v)$ of node v as the node w for which $r(\ell_w)$ is closest from above to node $r(\ell_v)$ (viewing $[0, 1)$ as a ring in both cases). Given two nodes v and w , we define their *distance* as

$$\delta(v, w) = \min\{(1 + r(\ell_v) - r(\ell_w)) \bmod 1, (1 + r(\ell_w) - r(\ell_v)) \bmod 1\}.$$

In order to maintain a cycle among the nodes, we simply have to maintain the following invariant:

Invariant 7.2 *Every node v in the system is connected to $\text{pred}(v)$ and $\text{succ}(v)$.*

Now, suppose that the labels of the nodes are generated via the recursive strategy above. Then we have the following properties:

Lemma 7.3 *Let n be the current number of nodes in the system, and let $\bar{n} = 2^{\lceil \log n \rceil}$. Then for every node $v \in V$:*

- $|\ell_v| \leq \lceil \log n \rceil$ and
- $\delta(v, \text{pred}(v)) \in [1/(2\bar{n}), 1/\bar{n}]$ and $\delta(v, \text{succ}(v)) \in [1/(2\bar{n}), 1/\bar{n}]$.

So the nodes are approximately evenly distributed in $[0, 1)$ and the number of bits for storing a label is almost as low as it can be without violating the uniqueness requirement. But how does the supervisor maintain the cycle? This is implied by the following demand, where n is again the current number of nodes in the system.

Invariant 7.4 *At any time, the supervisor stores the contact information of $\text{pred}(v)$, v , $\text{succ}(v)$, and $\text{succ}(\text{succ}(v))$ where v is the node with label $\ell(n - 1)$.*

In order to satisfy Invariants 7.2 and 7.4, the supervisor performs the following actions. If a new node w joins, then the supervisor

- informs w that $\ell(n)$ is its label, $\text{succ}(v)$ is its predecessor, and $\text{succ}(\text{succ}(v))$ is its successor,
- informs $\text{succ}(v)$ that w is its new successor,
- informs $\text{succ}(\text{succ}(v))$ that w is its new predecessor,
- asks $\text{succ}(\text{succ}(v))$ to send its successor information to the supervisor, and
- sets $n = n + 1$.

If an old node w leaves and reports ℓ_w , $\text{pred}(w)$, and $\text{succ}(w)$ to the supervisor (recall that we are assuming graceful departures), then the supervisor

- informs v (the node with label $\ell(n - 1)$) that ℓ_w is its new label, $\text{pred}(w)$ is its new predecessor, and $\text{succ}(w)$ is its new successor,

- informs $\text{pred}(w)$ that its new successor is v ,
- informs $\text{succ}(w)$ that its new predecessor is v ,
- informs $\text{pred}(v)$ that $\text{succ}(v)$ is its new successor,
- informs $\text{succ}(v)$ that $\text{pred}(v)$ is its new predecessor,
- asks $\text{pred}(v)$ to send its predecessor information to the supervisor and to ask $\text{pred}(\text{pred}(v))$ to send its predecessor information to the supervisor, and
- sets $n = n - 1$.

A detailed implementation of the leave and join operations can be found in Figures 1 and 2. In this implementation, we assume for simplicity that references to relay points can be freely exchanged, i.e., identities are not needed. It will be an assignment to implement the join and leave operations with the identity concept. The following lemma is not difficult to check and will also be an assignment.

Lemma 7.5 *The join and leave operations preserve Invariants 7.2 and 7.4.*

Hence, we arrive at the following theorem, which implies that our central requirements on a supervisor are kept.

Theorem 7.6 *At any time, the supervisor only needs to store the current value of n and a constant amount of contact information, and the join and leave operations only need a constant amount of messages and three communication rounds to complete.*

7.3 Concurrency

The above scheme only allows the supervisor to execute join and leave operations in a strictly sequential manner because it only has sufficient information to integrate or remove one peer at a time. In order to be able to handle d join or leave requests in parallel, we extend Invariant 7.2 with the following rule:

Invariant 7.7 *In addition to Invariant 7.2, every node v in the system is connected to its d th predecessor $\text{pred}_d(v)$ and its d th successor $\text{succ}_d(v)$.*

Furthermore, given that v is the node with label $\ell(n - 1)$, Invariant 7.4 needs to be extended to:

Invariant 7.8 *At any time, the supervisor stores the contact information of v , the $2d$ successors of v , and the $3d$ predecessors of v .*

These invariants are preserved in the following way.

Concurrent Join Operation. In the following, let v be the node with label $\ell(n - 1)$. Let $\text{succ}_i(v)$ denote the i th successor of v on the cycle and $\text{pred}_i(v)$ denote the i th predecessor of v on the cycle.

Let the d new peers be w_1, w_2, \dots, w_d . Then the supervisor integrates w_i between $\text{succ}_i(v)$ and $\text{succ}_{i+1}(v)$ for every $i \in \{1, \dots, d\}$. As is easy to check, this will violate Invariant 7.7 for the $2d$ closest successors of v and the $d - 2$ closest predecessors of v . But since the supervisor knows all of these nodes, it can directly inform them about the change. In order to repair Invariant 7.8, the supervisor will request information about the d th successor from the d furthest successors of v and will set v to w_d .

<pre> Supervisor { Supervisor() { n := 0 # counter v := NULL # node with label $\ell(n - 1)$ pv := NULL # pred(v) sv := NULL # succ(v) ssv := NULL # succ(succ(v)) } Join(w: Relay) { if (n = 0) { w ← setup(0, w, w) pv := w v := w sv := w ssv := w } else { w ← setup($\ell(n)$, sv, ssv) sv ← setSucc(w) ssv ← setPred(w) pv := sv v := w sv := ssv ssv := ssv ← getSucc() } n := n + 1 } } </pre>	<pre> Leave(ℓ: Int, pw: Relay, sw: Relay) { if (n > 0) { if (n = 1) { pv := NULL, v := NULL sv := NULL, ssv := NULL } else { # remove v from the system pv ← setSucc(sv) sv ← setPred(pv) if (pw = v) pv := pw if (sw = v) sv := sw # move v into position of w if (v ≠ w) { v ← setup(ℓ, pw, sw) pw ← setSucc(v) sw ← setPred(v) } # update pointers if (pv = w) pv := v if (sv = w) sv := v ssv := sv sv := pv v := pv ← getPred() pv := pv ← getPredPred() } n := n - 1 } } </pre>
--	--

Figure 1: Operations needed by the supervisor to maintain a cycle.

Concurrent Leave Operation. Let the d peers that want to leave the system be w_1, w_2, \dots, w_d . For simplicity, we assume that they are outside of the peers known to the supervisor and that they are not in the neighborhood of each other, but our strategy below can also be extended to these cases. The strategy of the supervisor is to replace w_i by $\text{pred}_{2(i-1)}(v)$ for every i . As is easy to check, this will violate Invariant 7.7 for the d closest successors of v and the $3d$ closest predecessors of v . But since the supervisor knows all of these nodes, it can directly inform them about the change. In order to repair Invariant 7.7, the supervisor will request information about the d th predecessor from the d furthest predecessors of v and their d th predecessors and will set v to $\text{pred}_{2d}(v)$.

The operations have the following performance.

Theorem 7.9 *The supervisor needs at most $O(d)$ work and $O(1)$ time (given that the work can be done in parallel) to process d join or leave requests.*

Peer {	setup($\ell : \text{Int}, p : \text{Relay}, s : \text{Relay}$) {
Peer() {	$label := \ell$
$label := 0$ # label of peer v	$pred := p$
$succ := \text{NULL}$ # succ(v)	$succ := s$
$pred := \text{NULL}$ # pred(v)	}
$sr := \text{new Relay}()$ # relay point of v	setSucc($w : \text{Relay}$) {
}	$succ := w$
	}
Join($s : \text{Relay}$) { # relay of supervisor	setPred($w : \text{Relay}$) {
if ($s \neq \text{NULL}$) {	$pred := w$
$s \leftarrow \text{Join}(sr)$	}
$super := s$ # current supervisor	getSucc(): Relay {
}	return $succ$
}	}
	getPred(): Relay {
Leave() {	return $pred$
if ($super \neq \text{NULL}$)	}
$super \leftarrow \text{Leave}(label, pred, succ)$	getPredPred(): Relay {
$super := \text{NULL}$	return $pred \leftarrow \text{getPred}()$
}	}

Figure 2: Operations needed by a peer to maintain a cycle.

7.4 Multiple Supervisors

If a supervised network becomes so large that a single supervisor cannot manage all of the join and leave requests, one can easily extend the supervised cycle to multiple supervisors. Suppose that we have k supervisors S_0, S_1, \dots, S_{k-1} . Then the $[0, 1)$ -ring is split into the k regions $R_i = [(i-1)/k, i/k)$, $1 \leq i \leq k$, and supervisor S_i is responsible for region R_i . Every supervisor manages its region as described for a single supervisor above, and the borders are maintained by communicating with the neighboring supervisors on the ring. The supervisors themselves form a completely interconnected network.

Each time a new node v wants to join the system via some supervisor S_i , S_i forwards it to a random supervisor to integrate v into the system. Each time a node v under some supervisor S_i wants to leave the system, S_i contacts a random supervisor (which may also be itself) to provide a replacement node. Using standard Chernoff bounds, we get:

Theorem 7.10 *Let n be the total number of nodes in the system. Then it holds for every $i \in \{1, \dots, k\}$ that the number nodes currently placed in R_i is in the range $n/k \pm O(\sqrt{(n/k) \log k + \log k})$, with high probability.*

Hence, if n is sufficiently large compared to k , then the multi-supervised cycle has basically the same properties as the single-supervised cycle above.

7.5 Recursively maintaining a tree

The cycle has a low degree but its diameter and expansion are very bad. The simplest way of achieving a low diameter is to use a tree. Thus, next we discuss how to recursively maintain a tree. As for the cycle, our basic approach will be to preserve Invariant 7.1. We will also preserve Invariant 7.2, though the edges implied by this Invariant will not be part of the tree. But they will tremendously simplify the task of maintaining a tree, as we will see. Altogether, the following connectivity information has to be preserved.

Invariant 7.11 *Every node v in the system with label $\ell_v = (\ell_1 \dots \ell_d)$ is connected to*

1. $\text{pred}(v)$ and $\text{succ}(v)$ (to form a cycle) and
2. the nodes with labels $(\ell_1 \dots \ell_{d-2}1)$, $(\ell_1 \dots \ell_{d-1}01)$, and $(\ell_1 \dots \ell_{d-1}11)$, if they exist (to form a tree).

Suppose that this invariant is kept at any time. Then the following lemma follows.

Lemma 7.12 *At any time, the n nodes (apart from node 0) form a binary tree of depth $\lceil \log n \rceil - 1$.*

Proof. Consider a binary tree with n nodes, and label the edge to the left child of any node “0” and to the right child of any node “1”. Let the label t_v of every node v in this tree be the sequence of edge labels when moving along the unique path from the root to v . Then every node v with label $(\ell_1 \dots \ell_d)$ is connected to the node with label $(\ell_1 \dots \ell_{d-1})$ (its parent), if it exists, and is also connected to the nodes with labels $(\ell_1 \dots \ell_d 0)$ and $(\ell_1 \dots \ell_d 1)$ (its children), if they exist. Defining t_v as ℓ_v (the label of v in our network) without the least significant bit, we see that Invariant 7.11(2) fulfills the connectivity requirements of a tree. Since it follows from Lemma 7.3 that every node has a label of size at most $\lceil \log n \rceil$, the depth of the tree can be at most $\lceil \log n \rceil - 1$ (when ignoring node 0). \square

Next we specify the connectivity information the supervisor needs in order to maintain the tree.

Invariant 7.13 *At any time, the supervisor stores the contact information of $\text{pred}(v)$, v , $\text{succ}(v)$, and $\text{succ}(\text{succ}(v))$ where v is the node with label $\ell(n)$.*

Hence, the supervisor does not need any further connectivity information beyond what it needs for the cycle. In order to satisfy Invariants 7.11 and 7.13, the supervisor performs the following actions. If a new node w joins, then the supervisor

- informs w that $\ell(n+1)$ is its label, $\text{succ}(v)$ is its predecessor, and $\text{succ}(\text{succ}(v))$ is its successor, and $\text{succ}(v)$ resp. $\text{succ}(\text{succ}(v))$ is its parent (depending on $\ell(n+1)$),
- informs $\text{succ}(v)$ that w is its new successor,
- informs $\text{succ}(\text{succ}(v))$ that w is its new predecessor,

- asks $\text{succ}(\text{succ}(v))$ to send its successor information to the supervisor, and
- sets $n = n + 1$.

Hence, from the point of view of the supervisor, the inclusion of a new node is almost identical to the cycle.

If an old node w leaves and reports ℓ_w , $\text{pred}(w)$, $\text{succ}(w)$, $\text{parent}(w)$, $\text{lchild}(w)$, and $\text{rchild}(w)$ to the supervisor, then the supervisor again executes almost the same steps as for the cycle.

When using the code for the supervisor given in Figure 3 and the code for the peers given in Figure 4, it is not difficult to prove the following lemma. Notice that for simplicity, we assume again that relay points can be freely exchanged.

<pre> Supervisor { Supervisor() { n := 0 # counter v := NULL # node with label $\ell(n)$ pv := NULL # pred(v) sv := NULL # succ(v) ssv := NULL # succ(succ(v)) } Join(w: Relay) { if (n = 0) { w ← setup(0, w, w, NULL, NULL, NULL) pv := w v := w sv := w ssv := w } else { if ($\ell(n) \& 2 = 0$) { w ← setup($\ell(n)$, sv, ssv, ssv, NULL, NULL) ssv ← setRightChild(w) } else { w ← setup($\ell(n)$, sv, ssv, sv, NULL, NULL) sv ← setLeftChild(w) } sv ← setSucc(w) ssv ← setPred(w) pv := sv v := w sv := ssv ssv := ssv ← getSucc() } n := n + 1 } } </pre>	<pre> Leave(ℓ: Int, pw: Relay, sw: Relay, fw, lcw, rcw: Relay) { if (n > 0) { if (n = 1) { pv := NULL, v := NULL sv := NULL, ssv := NULL } else { # remove v from tree if ($\ell(n-1) \& 2 = 0$) sv ← setRightChild(NULL) else pv ← setLeftChild(NULL) pv ← setSucc(sv) sv ← setPred(pv) if (pw = v) pw := pv if (sw = v) sw := sv if (lcw = v) lcw := NULL if (rcw = v) rcw := NULL # move v into position of w if (v ≠ w) { v ← setup(ℓ, pw, sw, fw, lcw, rcw) pv ← setSucc(v) sv ← setPred(v) if ($\ell \& 2 = 0$) fw ← setRightChild(v) else fw ← setLeftChild(v) if (lcw ≠ NULL) lcw ← setParent(v) if (rcw ≠ NULL) rcw ← setParent(v) } # update pointers if (pv = w) pv := v if (sv = w) sv := v ssv := sv sv := pv v := pv ← getPred() pv := pv ← getPredPred() } n := n - 1 } } </pre>
---	---

Figure 3: Operations needed by the supervisor to maintain a tree.

Lemma 7.14 *The join and leave operations preserve Invariants 7.11 and 7.13.*

Hence, we arrive at the following theorem.

<pre> Peer { Peer() { label := 0 # label of peer v succ := NULL # succ(v) pred := NULL # pred(v) parent := NULL lchild := NULL rchild := NULL sr := new Relay() # relay point of v } Join(s: Relay) { if (s ≠ NULL) { s → Join(sr) super := s # current supervisor } } Leave() { if (super ≠ NULL) super ← Leave(label, pred, succ, parent, lchild, rchild) super := NULL } setup(ℓ : Int, p : Relay, s : Relay, f: Relay, lc: Relay, rc: Relay) { label := ℓ pred := p succ := s parent := f lchild := lc rchild := rc } </pre>	<pre> setSucc(w: Relay) { succ := w } setPred(w: Relay) { pred := w } setParent(w: Relay) { parent := w } setLeftChild(w: Relay) { lchild := w } setRightChild(w: Relay) { rchild := w } getSucc(): Relay { return succ } getPred(): Relay { return pred } getPredPred(): Relay { return pred ← getPred() } </pre>
---	---

Figure 4: Operations needed by a peer to maintain a tree.

Theorem 7.15 *At any time, the supervisor only needs to store the current value of n and a constant amount of contact information, and the join and leave operations only need a constant amount of messages and three communication rounds to complete.*

Broadcasting

The dynamic tree can be used for efficient broadcasting. Suppose that some node v wants to broadcast information to all other nodes in the system. One way of solving this is that it forwards the broadcast message directly to the supervisor (so that the supervisor can authorize the broadcast, for example) and the supervisor initiates sending the broadcast message down the tree. A prerequisite for this is that the supervisor remembers the node with label 0, called *root* by it. If this is the case, then the code in Figure 5 will be executed correctly.

Inspecting the code, we arrive at the following result, which is optimal for broadcasting in constant

```

# operations of supervisor

Broadcast( $m$  : Message) {
     $root \leftarrow \text{sendDown}(m)$ 
}

# operations of peer

Broadcast( $m$  : Message) {
    if ( $super \neq \text{NULL}$ )  $super \leftarrow \text{Broadcast}(m)$ 
}

sendDown( $m$  : Message) {
    if ( $lchild \neq \text{NULL}$ )  $lchild \leftarrow \text{sendDown}(m)$ 
    if ( $rchild \neq \text{NULL}$ )  $rchild \leftarrow \text{sendDown}(m)$ 
    # handle broadcast message
}

```

Figure 5: Implementation of a broadcast operation in the dynamic tree.

degree networks. Here, the *dilation* means the longest path taken by a message in the broadcast operation.

Theorem 7.16 *The broadcast operation has a dilation of $O(\log n)$ and requires a work of $O(n)$.*

References

- [1] K. Kothapalli and C. Scheideler. Supervised peer-to-peer systems. In *Proc. of the 2005 International Symposium on Parallel Architectures, Algorithms, and Networks (ISPAAN)*, 2005.
- [2] C. Riley and C. Scheideler. A distributed hash table for computational grids. In *18th Int. Parallel and Distributed Processing Symposium (IPDPS)*, 2004.
- [3] C. Riley and C. Scheideler. Guaranteed broadcasting using SPON: A supervised peer overlay network. In *3rd International Zürich Seminar on Communications (IZS)*, 2004.