

# Distributed Systems: Distributed Hash Tables

**Gian Pietro Picco**

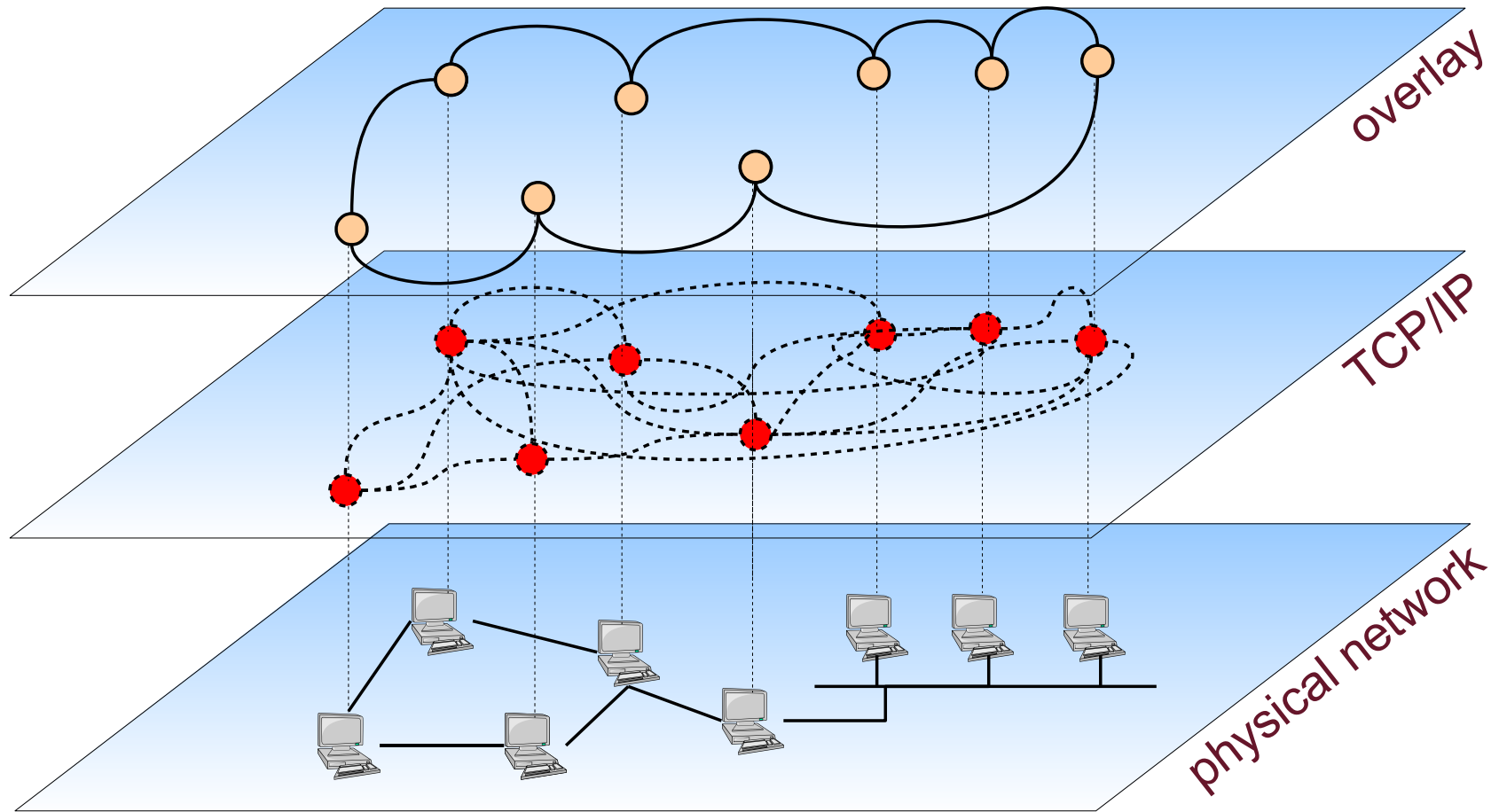
Dipartimento di Ingegneria e Scienza dell'Informazione  
University of Trento, Italy

`gianpietro.picco@unitn.it`  
`http://disi.unitn.it/~picco`

# Distributed Hash Tables

- A hash table is a data structure that efficiently associates a **key** to an **item**
  - `put(item, id)`
  - `item = get(id)`
  - A key is generated by applying a **hash function** to an item (e.g., SHA-1)
    - The mathematical properties of the function guarantee that it is very unlikely (although possible) that different items are hashed to the same key
- In a **distributed hash table (DHT)** the key and items are distributed across network hosts
  - “items” are the identifiers of resources
    - e.g., files, users
  - “keys” are associated to resources **and** to hosts
  - the name space is flat: only “opaque” identifiers, no structure in the names to be searched (unlike DNS)

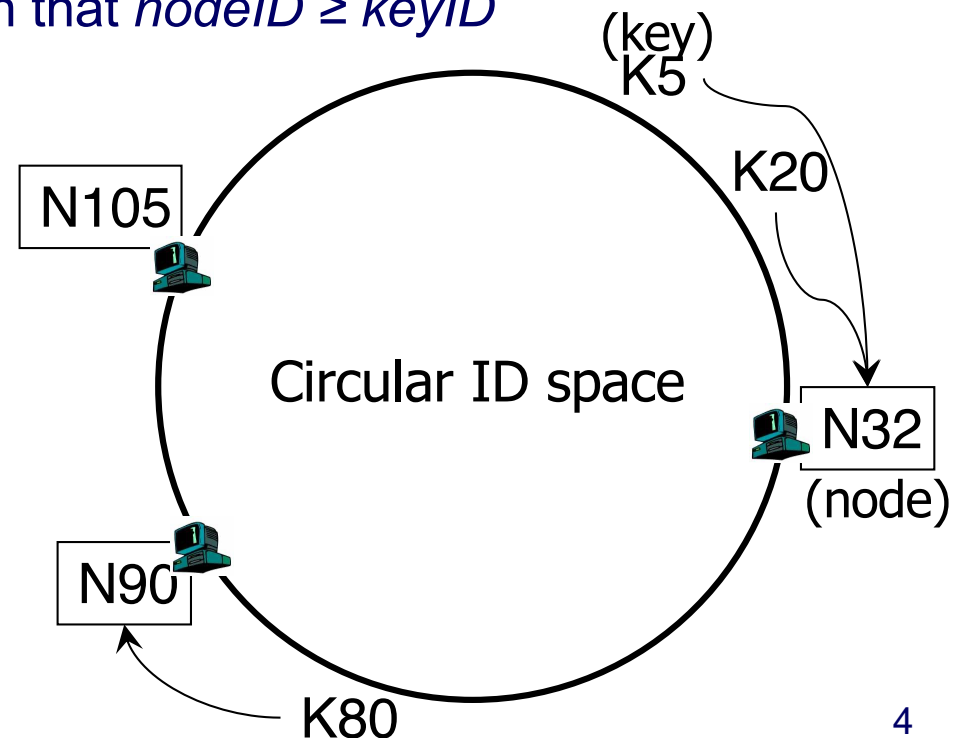
# Overlay Networks



- Tremendous design flexibility in providing application-specific topologies and protocols

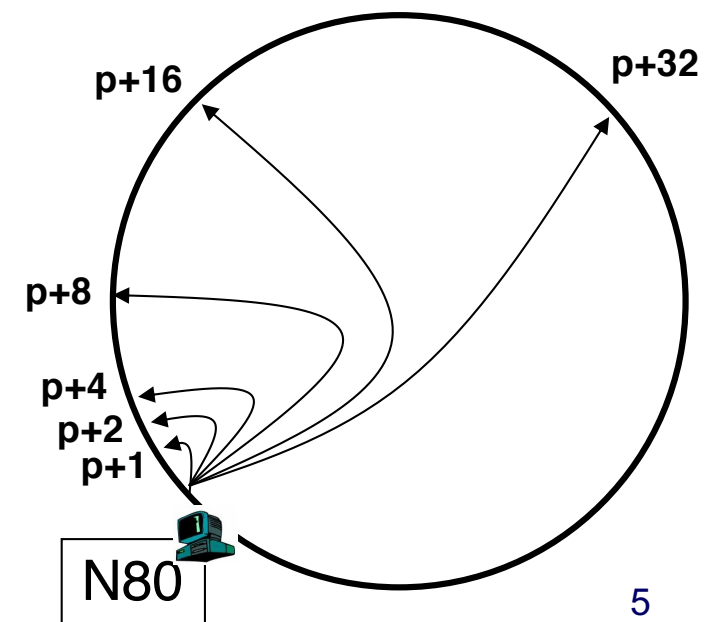
# DHT Example: Chord

- Idea: associate to each node and item a unique ID in a uni-dimensional space (a ring)
  - e.g., in the range  $[0...2^m]$ ,  $m$  is typically 160 bits
  - the ID is the hash of
    - the content of the item (e.g., a file)
    - the IP address, in the case of a node
  - A key  $k$  is **mapped** to its **successor**  $\text{succ}(k)$ , the node with next higher ID
    - i.e. the smallest node such that  $\text{nodeID} \geq \text{keyID}$
- Properties:
  - routing table size is  $O(\log N)$ , where  $N$  is #nodes
  - a file is found in  $O(\log N)$  hops

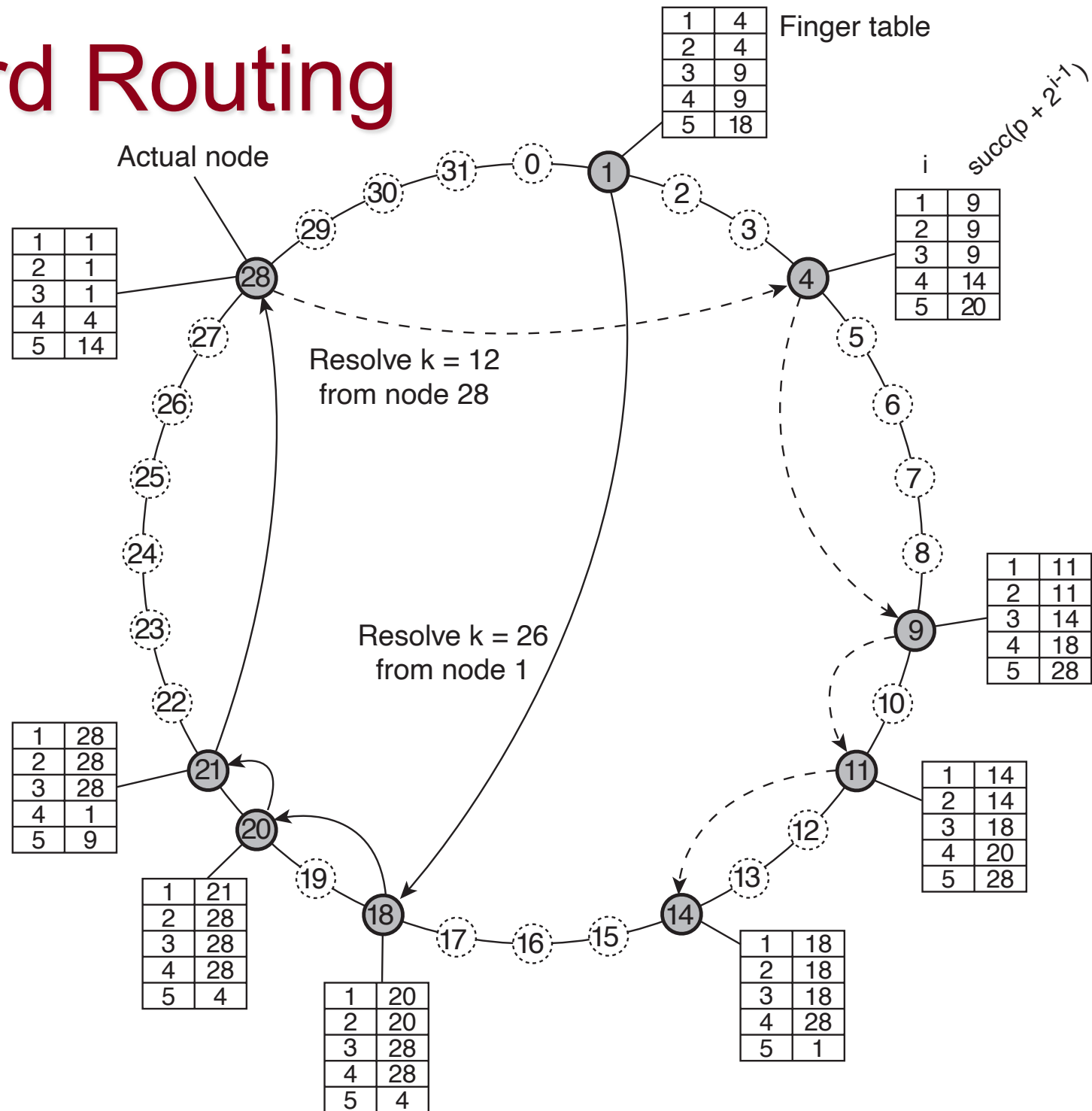


# Finger tables

- The difficult problem in DHTs is to efficiently search for the node hosting an item, given the item's key
  - $\text{lookup}(k)$  returns  $\text{succ}(k)$
  - Trivial, inefficient solution: linear search
    - each node  $p$  keeps track of  $\text{succ}(p+1)$  and  $\text{pred}(p)$
    - when  $p$  receives a  $\text{lookup}(k)$ , it serves it if  $\text{pred}(p) < k \leq p$ , otherwise it routes to  $\text{succ}(p+1)$  or  $\text{pred}(p)$
- In Chord, the problem is solved by using routing tables (“finger tables”) that contain “shortcuts” to exponentially bigger portions of the ID space
  - $\text{FT}_p(i) = \text{succ}(p+2^{i-1})$ 
    - Entry  $i$  in the finger table  $\text{FT}_p$  of node  $p$  is the first node that succeeds  $p$  by at least  $2^{i-1}$
    - Modulo  $m$ , restarts after 0
  - A  $\text{lookup}(k)$  is routed to the node  $q$  with index  $j$  in the FT such that  $\text{FT}(j) \leq k < \text{FT}(j+1)$



# Chord Routing





# Dealing with Churn

- Churn is the rate at which nodes join and leave the system
  - This may be voluntary (e.g., user disconnects) or induced by (communication or host) failures
- If a node  $p$  wants to join, it requests a lookup for  $\text{succ}(p+1)$ ; then it contacts this node and its predecessor to insert itself in the ring
  - Similar for leaving
- Problem: keep the finger tables up to date
  - Each node runs a background procedure that contacts  $\text{succ}(q+1)$  and asks to return  $\text{pred}(\text{succ}(q+1))$ 
    - $\text{FT}_q[1]$  must refer to  $\text{succ}(q+1)$
    - If the latter returns  $q$ , everything is consistent
    - If the successor returns  $p$  such that  $q < p \leq \text{succ}(q+1)$  a node  $p$  has joined;  $q$  updates its routing table to set  $\text{FT}_q[1]=p$
  - Similar procedures are run in the background to verify the consistency of the other finger tables, as well as  $\text{pred}(q)$