

# Peer-to-Peer Systems

Winter semester 2014

Jun.-Prof. Dr.-Ing. Kalman Graffi

Heinrich Heine University Düsseldorf

# Peer-to-Peer Systems

## Structured P2P Overlays

- Kademlia

## Interconnection Overlay Networks

- Hypercube
- De Bruijn Network

# Peer-to-Peer Systems

## Structured Homogenous P2P Overlay Networks – Kademlia

## “Kademlia” vs. “Kad”:

- Kademlia: Algorithm presented by Maymounkov and Mazières
- Kad: Implementation of the Kademlia algorithm used by eMule
- BitTorrent clients also implement Kademlia for their “trackerless mode”
  - Not compatible with Kad

## Uses 160-bit-IDs

## Uses XOR metric to determine distance between IDs

- XOR = bitwise eXclusive OR. Example:
- 10101010 = 170
- XOR     10010101 = 149
- =         00111111 = 63

## XOR metric is symmetric: $A \text{ XOR } B == B \text{ XOR } A$

- Contrast to Chord

The Kademlia protocol consists of 4 RPCs:

- FIND\_NODE(KEY):
  - Recipient returns <IP Address, UDP Port, Node ID> triples for k closest nodes he knows about
  
- FIND\_VALUE(KEY):
  - Like FIND\_NODE
  - With exception:
    - If recipient already stores the value, the value is returned instead of k closest nodes
  
- PING(IP ADDRESS)
  - Probes a node to see if it is online
  
- STORE(KEY, VALUE)
  - Instructs a node to store a <key,value> pair

# Node state

---

Each node maintains routing table (k buckets)

- Routing tables of different peers may be different

For each  $0 \leq i < 160$  every node

- keeps a list of <IP Address, UDP Port, Node ID> triples
- for k nodes within distance  $[2^i ; 2^{(i+1)}[ \rightarrow$  routing bucket I
  - Distance of  $2^i$  is the same as common prefix of i-1 bits
- in total  $k * 160$  contacts
- often: tree presentation of routing buckets / other prefix-subtrees

Nodes learn from

- messages they receive or
- using the FIND\_NODE method

Preference towards old contacts

- Study has shown that the longer a node has been up, the more likely it is to remain up another hour
- Resistance against DoS attacks by flooding the network with new nodes

# k Buckets: Peer Selection Policy

---

If node  $u$  receives a message from node  $v$   
then it adds node  $v$  to its  $k$ -bucket  
according to the following rules:

- IF  $v$  is already in a  $k$ -bucket  
THEN move  $v$  to the tail of the bucket
- IF  $v$  is not in the  $k$ -bucket and the bucket has fewer than  $k$  entries  
THEN insert recipient to the tail of list
- IF the appropriate  $k$  bucket is full AND least recently seen node is alive  
THEN move least recently seen node to tail of bucket and put node  $v$   
to waiting list for the bucket
- IF the appropriate  $k$  bucket is full AND least recently seen node is not alive  
THEN remove least recently seen node from bucket and add node  $v$  at the tail

Note: approx.  $k = 20$  in the real world

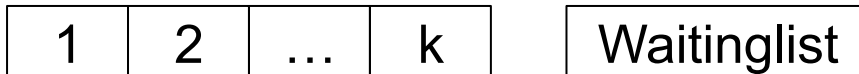
- Connection to nodes in buckets are NOT maintained

# Presentation of k buckets

## One presentation

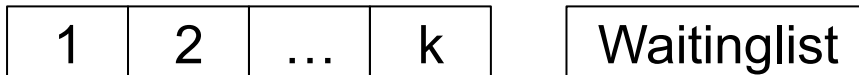
Bucket for  $2^{159} - 2^{160}$

Shared prefix = 0 bits



Bucket for  $2^{158} - 2^{159}$

Shared prefix = 1 bit



...

Bucket for  $2^0 - 2^1$

Shared prefix = 159 bits





The lookup initiator sends a FIND\_NODE request to  $\alpha$  nodes closest to target ID

- $\alpha$  (alpha) is the parameter for the number of parallel lookups
- Parallel asynchronous queries
  - Overcome faulty nodes

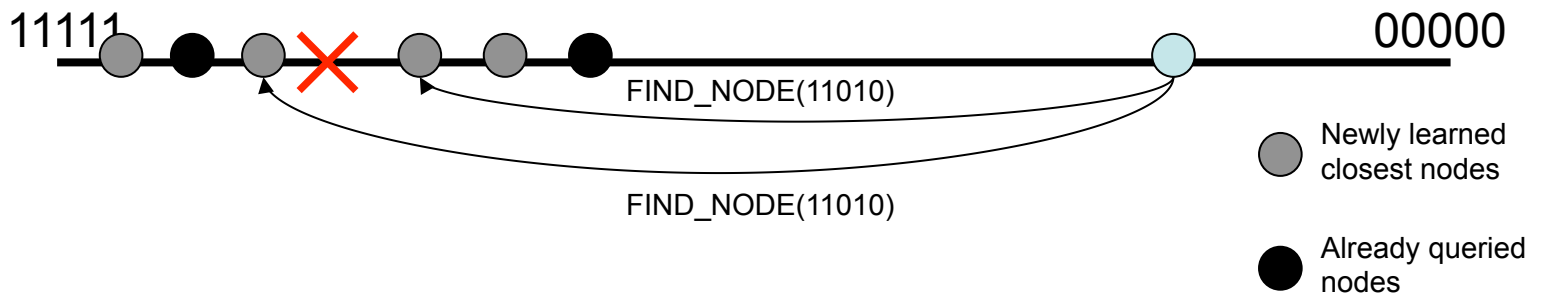
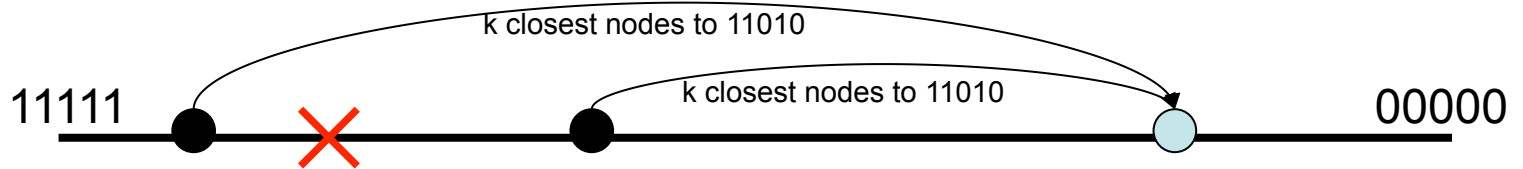
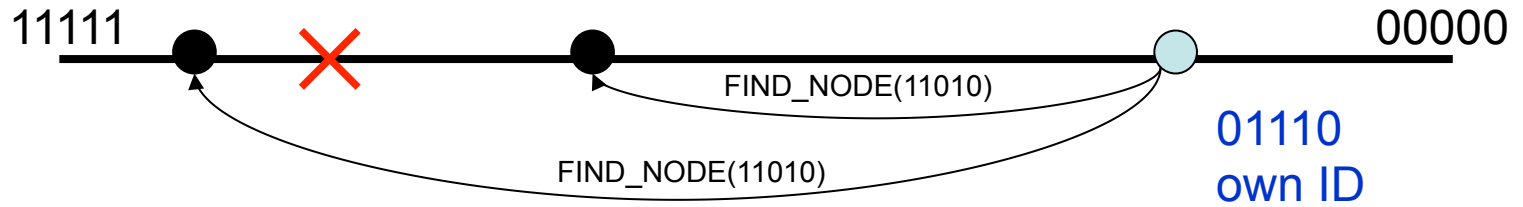
Receiving nodes respond with  $k$  nodes closest to target ID

Lookup initiator resends FIND\_NODE request to  $\alpha$  nodes learned from previous requests

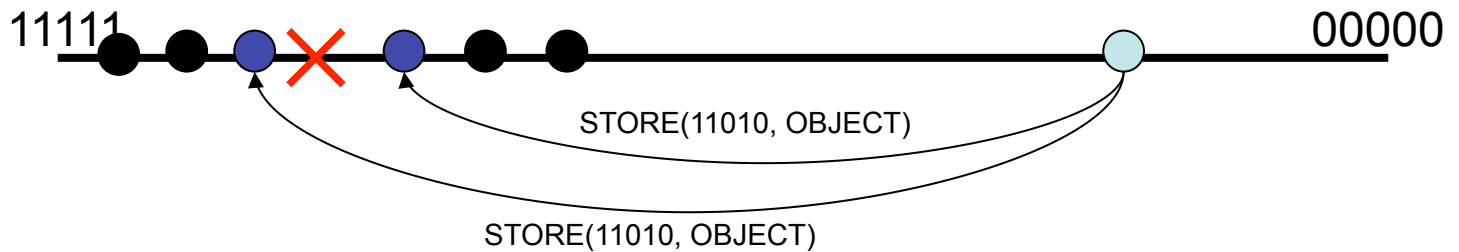
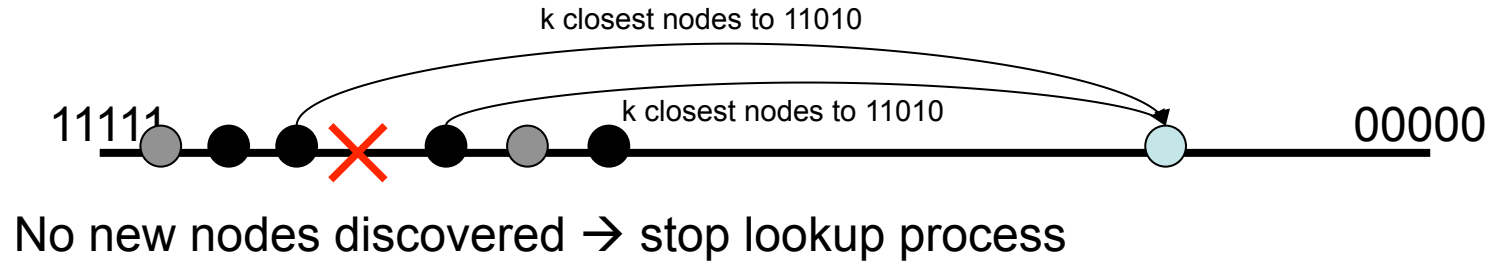
Procedure stops

if the set of  $k$  closest nodes does not change anymore

# Routing from 00010 to 11010, $\alpha = 2$ , $k = 2$



# Routing from 00010 to 11010, $\alpha = 2$ , $k = 2$



- Newly learned closest nodes
- Already queried nodes
- k closest nodes

## Parallel queries

- For one query,  $\alpha$  (alpha) concurrent lookups are sent
- More traffic load, but lower response times

## Network maintenance

- In Chord: active fixing of fingers
- In Kademlia: learning for bypassing queries
- Check if peer IDs fit better in routing table

## Large routing tables

- In Chord: 1 finger per distance  $2^i$  to  $2^{(i+1)}$
- In Kademlia:  $k$  contacts per distance  $2^i$  to  $2^{(i+1)}$
- Increased robustness

# Kademlia Limitations

## Kademlia functions

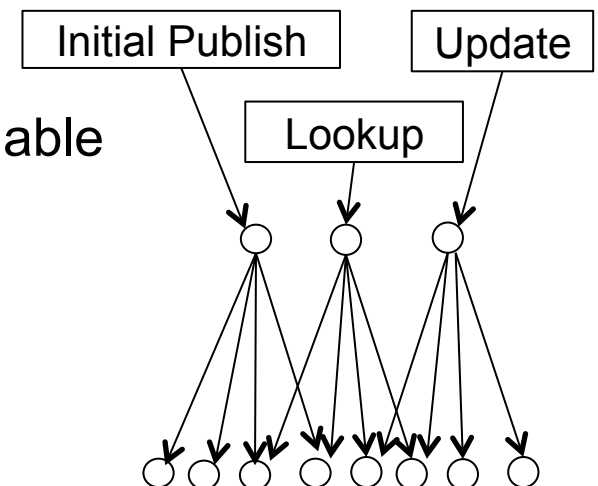
- Publish (key, value) pair ONCE
- Retrieve value for given key

## Different than goals than Pastry / Key-based Route

- Pastry: Find (single) responsible peer for given key
- Pastry: Allows update of data!

## Kademlia limitations

- Only 18% of nodes storing key/value reachable
- No updates possible (replication difficult)
- Key/values to be refreshed every 24h



# Kad – Introduction

---

Kad is the Kademlia implementation of the file sharing applications  
eMule/aMule

Uses UDP for communication

Has been developed ...

- ... to improve source lookup and
- ... to become independent from servers

Kad can be used ...

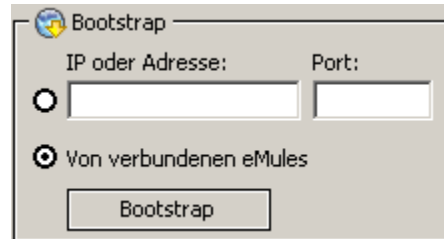
- ... to find sources for files and
- ... to perform keyword searches

Differences to Kademlia

- ID length
  - 128 bits due to MD4 hashing of eDonkey network (usually displayed in hex)
- Lookup
  - Two phases (1. Node lookup; 2. Action, e.g. source request)
  - Does not terminate when one result is found
- Variable k:  $k = 10$  ( $\alpha$  is still 3)

## Bootstrap

- A node can bootstrap using ...
  - ... a node known via the eDonkey network
  - ... a user-provided host that is already part of the Kad network
  - ... a downloaded hosts.dat file



# Kad – Publication and search

---

## Source publication

- A client obtains IDs for a file by hashing
  - The file content → 1st level pointer
  - Each part of the file name → 2nd level pointer
- And announces the file at peers responsible for the IDs

## Source search

- If ID of a file is known
- Kad can be used to search for peers that offer that file

## Keyword search

- Truncate search term at space characters → string tokens
- Hash tokens, lookup for token ID → obtain fitting file names
- Return (IDs, filenames)
- User chooses at GUI which files to download



# Kad – Keyword search example

The screenshot shows the eMule v0.48a search interface. The search term 'ubuntu' is entered in the search bar. The search results are displayed in a table with columns for filename, size, availability, completeness, type, file ID, and artist.

Dateiname	Größe	Verfü...	Vollst...	Typ	Datei ID	Künstl...
Ubuntu 7.10 i386.iso	4.23 GB	1	?	CD-Images	8D19D55D9C23C9BD9ABBA0075DBA3D56	
ubuntu-7.04-dvd-i386.iso	3.97 GB	32	?	CD-Images	1D297CF1AC2359B7FA4D432C4B164D0A	
[Ubuntu.Feisty.Fawn.7.04].ubuntu-7.04-...	3.96 GB	14	?	CD-Images	978E6D2137BCB087E2A96FC4FF6722D6	
Video2Brain - Ubuntu 6.06 Schulungs DVD...	3.68 GB	6	?	CD-Images	B7FB1B870A54FB76EDC8AEFAD1FAFB5C	
[Ubuntu.Linux.6.10.DVD].ubuntu-6.10-...	3.52 GB	4	?	CD-Images	2EFF4347F8A2A5DE76618D6FE270203A	
ubuntu-8.04.vdi	3.21 GB	1	?		9186265FA3B33CE8CE1BB2E7AAD1257B	

# Kad – Lookup procedure

---

At the beginning of a lookup procedure

- Searching node creates a search list
  - containing the  $k$  (= 10!) closest nodes to the target ID it knows about
- The list is ordered by XOR distance (closest at the top)

Lookup procedure (repeat until closest  $k$  nodes do not change):

- Node picks  $\alpha$  (= 3) random nodes from the search list
- Searching node sends a REQUEST message to them
  - Asks for closer nodes
- The queried nodes answer
  - with nodes they know about which are closer to the target than themselves
- Those nodes are inserted into the search list.

When the closest  $k$  nodes have been found:

- the specific action is performed
- (e.g. SEARCH REQUEST or PUBLISH REQUEST)



## Advantages

- Has successfully proven function in big real world scenarios
- Resistant to node failures
- Proximity of nodes is used

## Drawbacks

- Large routing table required
- Not deterministic routing:
  - Node A performs a lookup: finds k closest nodes
  - Node B performs a lookup: finds “other” k closest nodes
    - k closest node not necessarily the SAME!

## Future work required

- Security issues
  - Node IDs can be chosen freely

# Peer-to-Peer Systems

## Interconnection Overlay Networks – Hypercube Networks

## Motivation

- Assumption 1: no churn
- Assumption 2: number of nodes controllable
- Goal: overlay network for ID-based routing
- Efficient network topologies are possible

## Efficient topologies

- Fixed node degree
- Routing in  $O(\log N)$

## Differences to DHTs

- DHT: provides key-value lookup functionality
- Interconnection networks: only ID-based routing

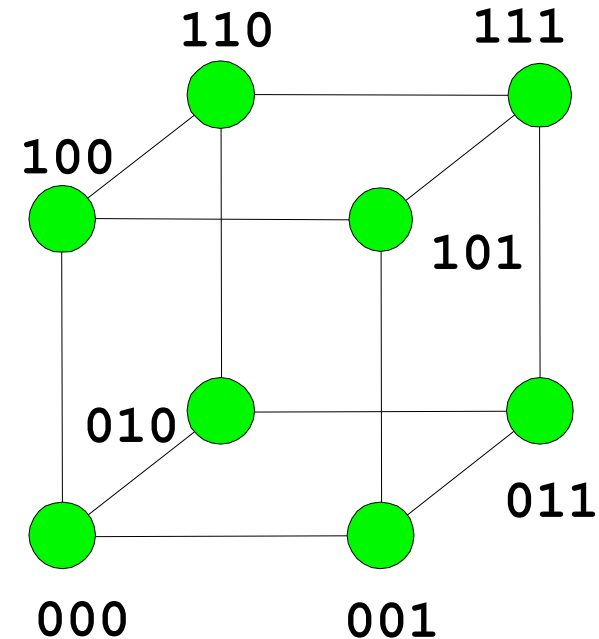
# Hypercube Networks

## Example:

- 3-dimensional Hypercube Network

## d-dimensional binary hypercube definition

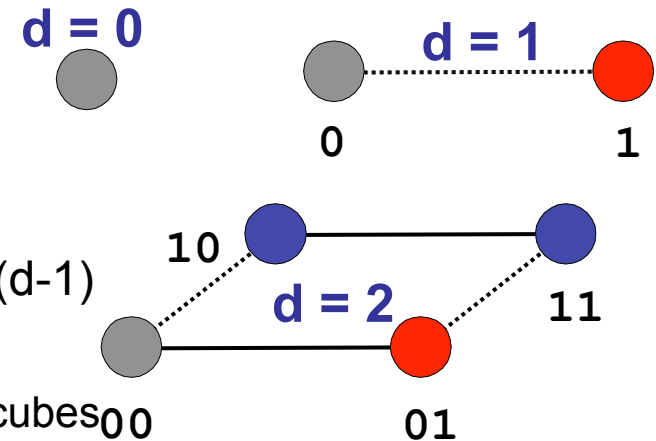
- A d-dimensional binary hypercube network consists of  $2^d$ 
  - interconnected nodes
  - whose addresses are represented by d-bit binary numbers.
- A node n has d adjacent nodes
  - whose addresses are obtained by reverting each bit of its address
  - fixed number of contacts per node!



# Hypercube Networks: How to construct

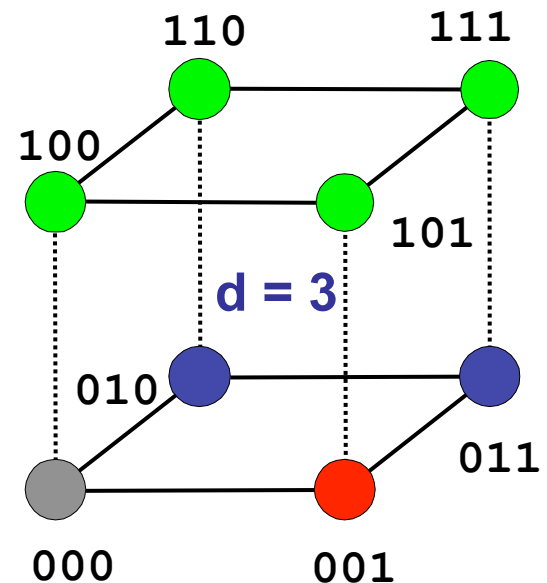
## Construction algorithm

- if ( $d == 0$ )
  - draw a node
- else
  - 1. draw two hypercubes of dimension ( $d-1$ )
  - 2. add arcs between
    - corresponding nodes of the two hypercubes



## Identification algorithm

- 1. use the first bit
  - to decide which of the two sub-hypercubes the node is in.
- 2. use the remaining  $d-1$  bits as
  - an address within that sub-hypercube.





A simple source-to-destination routing algorithm for binary hypercube networks

- Node failure is ignored

To route from node  $i$  to node  $j$ ,

- the bits of  $i$  which are different from those of  $j$  are changed,
  - traversing a link with each bit changed,
  - until  $j$  is reached.
- Since each node has links corresponding to each of the  $n$  bit positions,
  - any bit can be changed at any step in the route;
- →
  - the number of eligible links at each step equals the distance from the destination.

# Hypercube Networks: Routing Example

## Route from 010 to 101

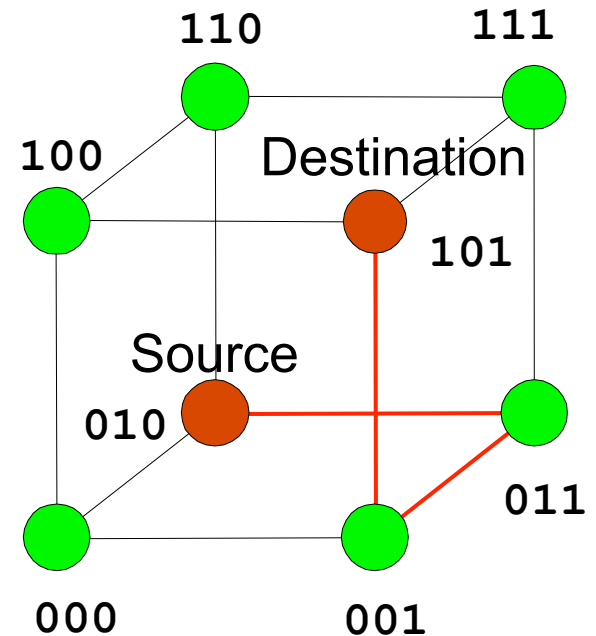
- 010 has 3 options (000, 110, 011)
  - → distance = 3
  - Select 011 as the next hop
- 011 has 2 options (001, 111)
  - → distance = 2
  - Select 001 as the next hop
- 001 is direct neighbor
  - → distance = 1
  - Forward to destination

IF

- nodes fail

THEN

- alternative paths may be followed in dynamic environments



# Hypercube Networks: Properties

Basic properties:

Diameter of the network

- i.e. worst case node distance
- increases logarithmically with respect to the network size

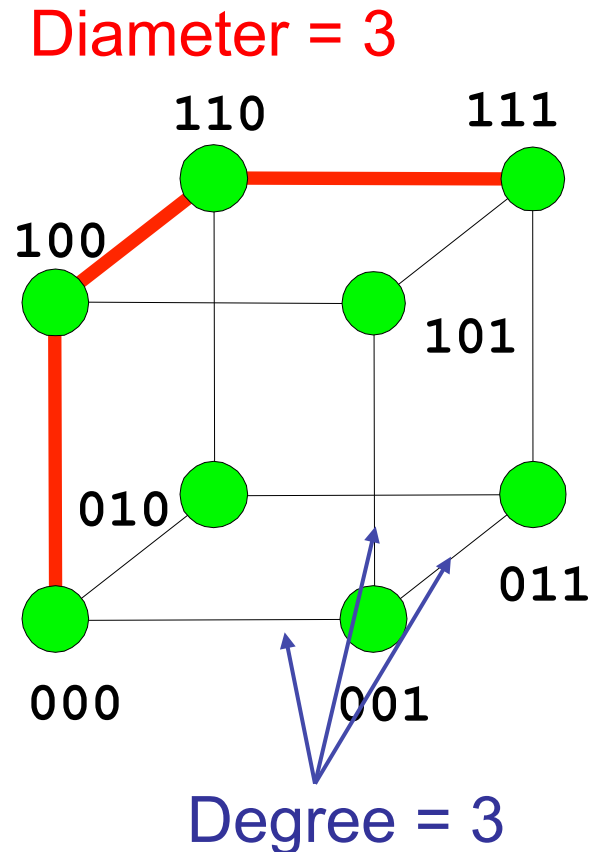
Node degree

- i.e. number of neighbors
- increases logarithmically with respect to the network size

Network is both

- vertex- and
- edge-symmetric graph

Hypercube is a hierarchically recursive network



# Hypercube Networks: Limitations

---

## Exponentially expandable

- Network size is defined only for 1, 2, 4, 8, ... nodes
  - for binary hypercubes
- Not incrementally expandable
  - i.e. network cannot be defined for any arbitrary integer

## Node degree increases logarithmically with respect to the network size

- → increasing the required maintenance cost

## Mostly appropriate for static environments

- with infrequent joins and leaves

## Scenarios:

- Controlled environments (provider defined topologies)

## Benefits

- Fixed and low node degree
- Fast routing:  $O(\log N)$

## Drawbacks

- Exponentially expandable
  - Network size is defined only for binary graphs
    - Node count: 1, 2, 4, 8, 16 ...  $2^i$
  - Not incrementally expandable
    - i.e. basic network cannot be defined for any arbitrary integer
- Cannot handle churn

# Peer-to-Peer Systems

## Interconnection Overlay Networks – de Bruijn Networks

## Motivation

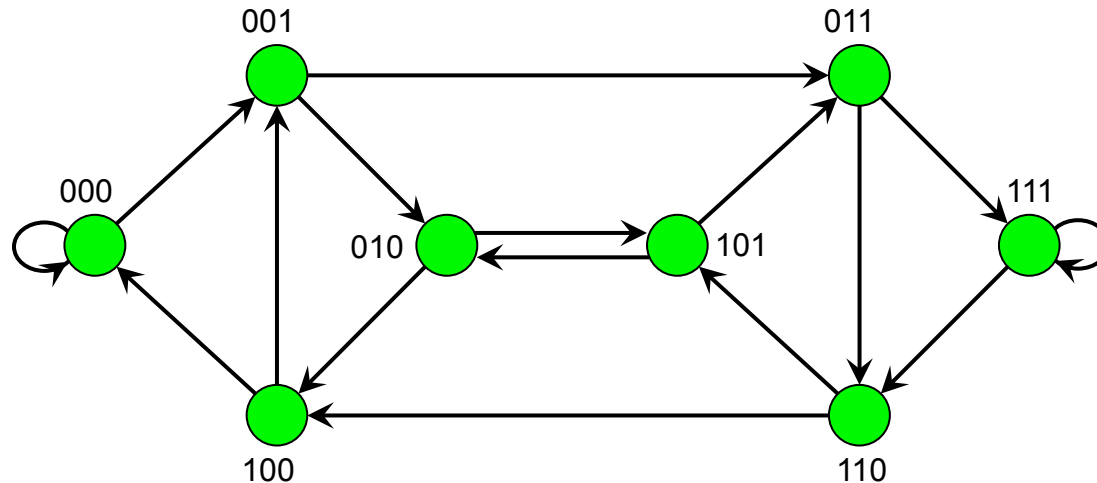
- Assumption 1: no churn
- Assumption 2: number of nodes controllable
- Goal: overlay network for ID-based routing
- Efficient network topologies are possible

## Efficient topologies

- Fixed node degree
- Routing in  $O(\log N)$

## Differences to DHTs

- DHT: provides key-value lookup functionality
- Interconnection networks: only ID-based routing



r-dimensional binary de Bruijn digraph (directed graph) consists of

- $2^r$  nodes and  $2^{r+1}$  edges

Each node corresponds to

- An r-digit binary string

There is a directed edge

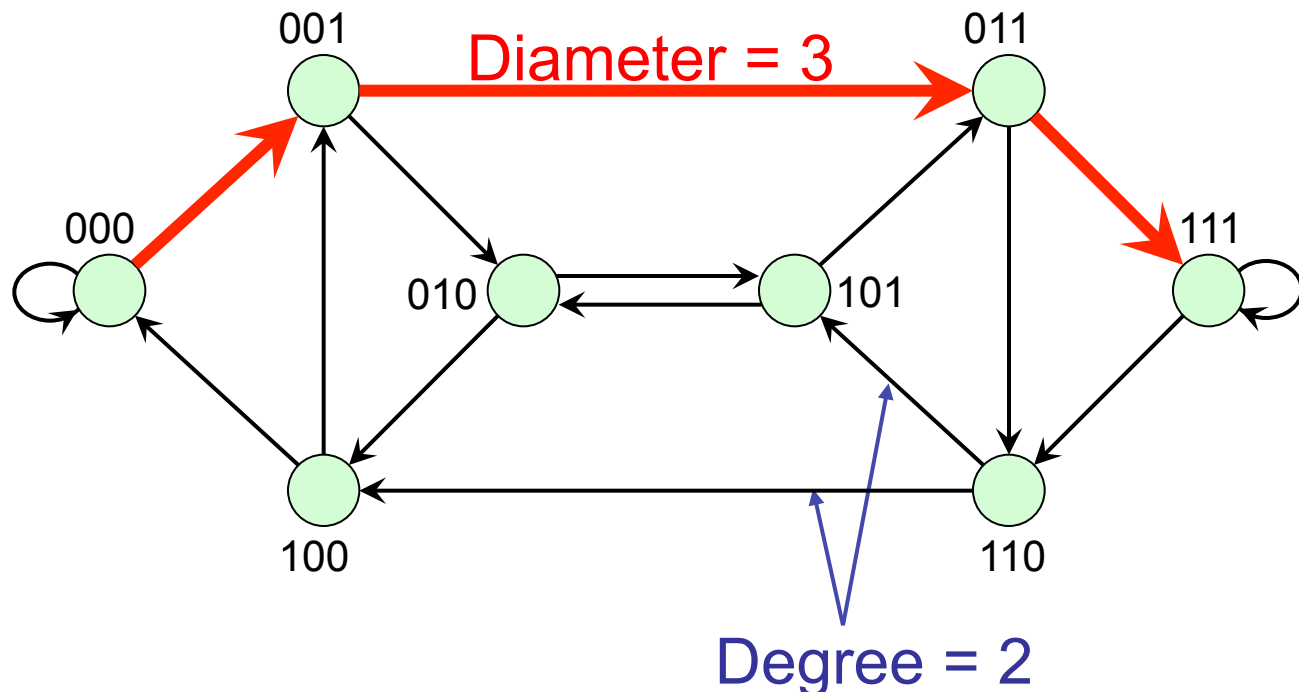
- from each node  $(u_1 u_2 \dots u_r)$  with  $u_i$  either 0 or 1
- To  $(u_2 u_3 \dots u_{r-1} 0)$  and  $(u_2 u_3 \dots u_{r-1} 1)$
- Principle: drop first character and add 0 or 1 to the end

Number of in and out edges per node: 2 each



## Basic characteristics

- Average distance is very close to the diameter
- Constant vertex degree
- Logarithmic diameter
- Adjacency is based on left shift by 1 position



# de Bruijn Networks: Routing Algorithm

## Assumptions

- Dimension of the network
- Source node  $K$
- Destination node  $L$

Operation:  $\text{shift\_match}(s, K, L)$ , where  $0 \leq \text{shift} \leq d$

- returns TRUE if and only if  $k_{1+\text{shift}}k_{2+\text{shift}}\dots k_D = l_1l_2\dots l_{D-\text{shift}}$
- returns FALSE otherwise
  - ~ colloq: cut the same no. of bits in the source ID left as in destination ID right

Operation:  $\text{merge}(\text{shift}, K, L)$ , where  $0 \leq \text{shift} \leq D$ .

- returns the sequence of length  $(D + \text{shift})$  given by  $k_1k_2\dots k_Dl_{D-\text{shift}+1}l_{D-\text{shift}+2}\dots l_D$

Shortest-path algorithm (Sivarajan and Ramaswami)

$\text{shortest\_path}(K, L)$

- $\text{shift} = 0$
- while ( $\text{shift\_match}(\text{shift}, K, L) == \text{FALSE}$  and  $\text{shift} < D$ )
- do  $\text{shift} = \text{shift} + 1$
- return  $\text{merge}(\text{shift}, K, L)$

# de Bruijn Networks: Routing Example

Example: route

- from node with ID K: 001
- to node L: 101

Common match: 00\_\_1\_\_01

Only shift left operation is allowed

- i.e. only 2 entries in the routing table

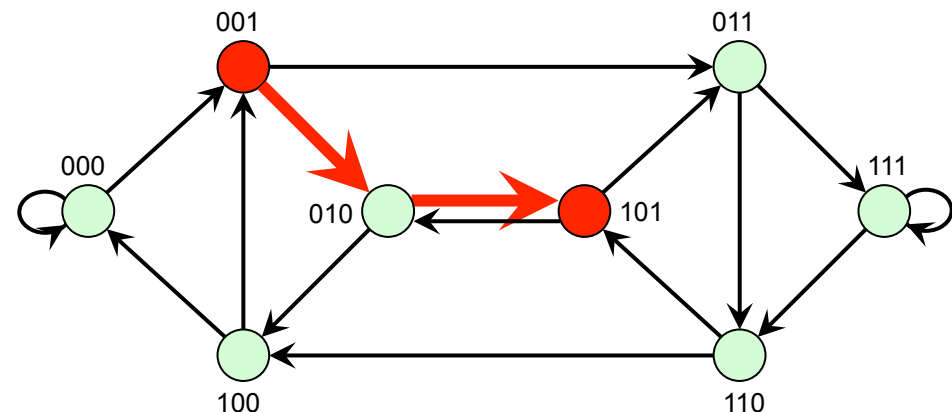
Node address			Link
0	0	1	0
0	1	0	1
1	0	1	

Example

from node 001  
to node 101

Left – shift operation:

- Drop left letter
- Add new letter on right side
  - Corresponds to using link



# de Bruijn Networks: Routing Example

i.e. Route message from source (001) to (101)

shift = 0

shift\_match(0, 001, 101) = FALSE ... 001 101

shift = 1

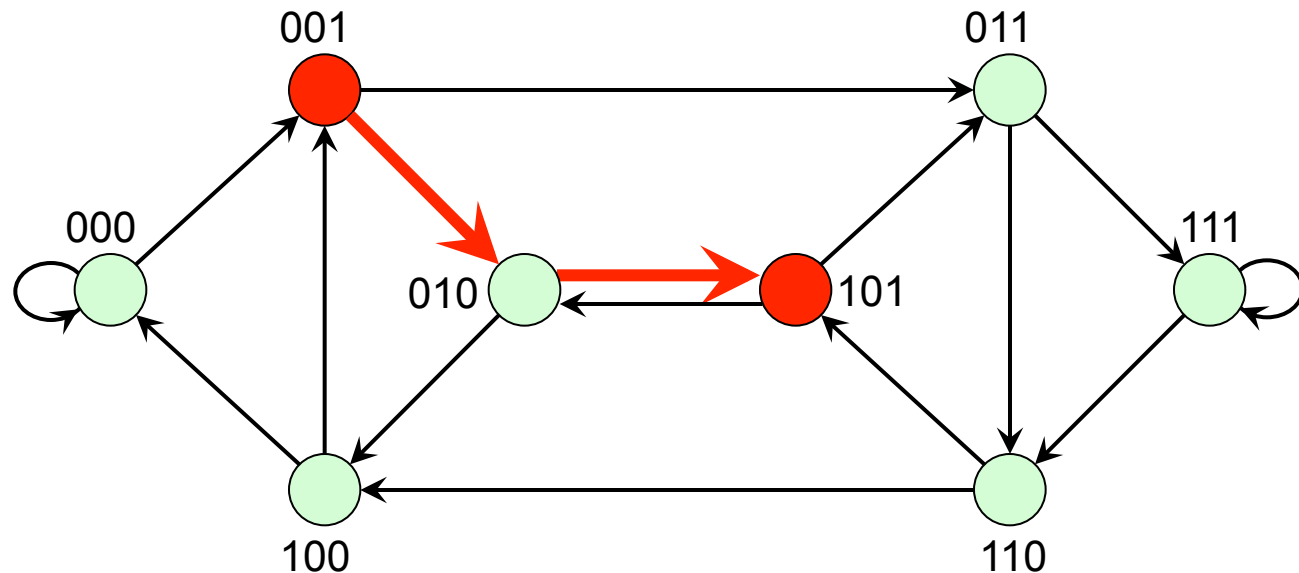
shift\_match(1, 001, 101) = FALSE ... X01 10X

shift = 2

shift\_match(2, 001, 101) = TRUE ... XX1 1XX

return merge (2, 001,101) = 00101 ... 01

- 2 forwarding steps



# de Bruijn Networks: Routing Example

## Example: route

- from node with ID 11110
- to node 01001

## Common match:

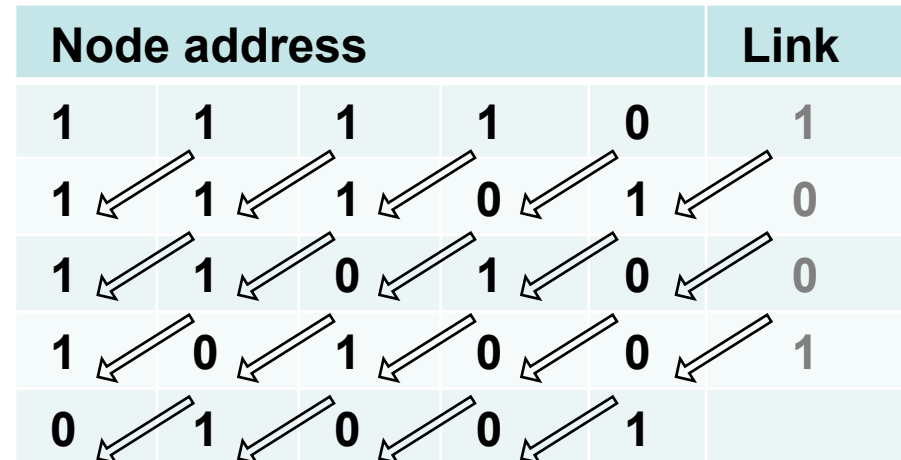
1111\_0\_1001

Bits to insert stepwise: 1001

Only shift left operation is allowed

- i.e. only 2 entries in the routing table

Node address					Link
1	1	1	1	0	1
1	1	1	0	1	0
1	1	0	1	0	0
1	0	1	0	0	1
0	1	0	0	1	



## Partial Line Digraph Algorithm: recursive construction

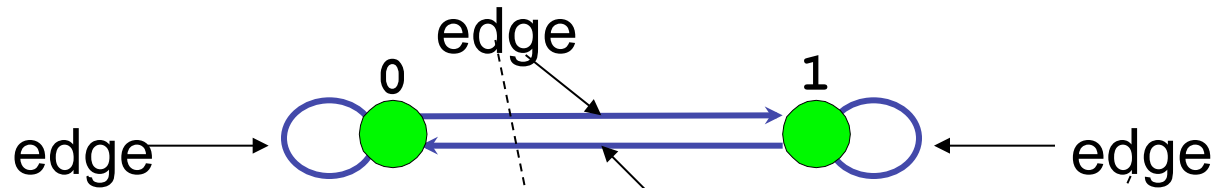
- The  $N$ -node graph (dimension  $r+1$ ) can be obtained from two  $N/2$ -node graphs (dimension  $r$ )

## Algorithm description

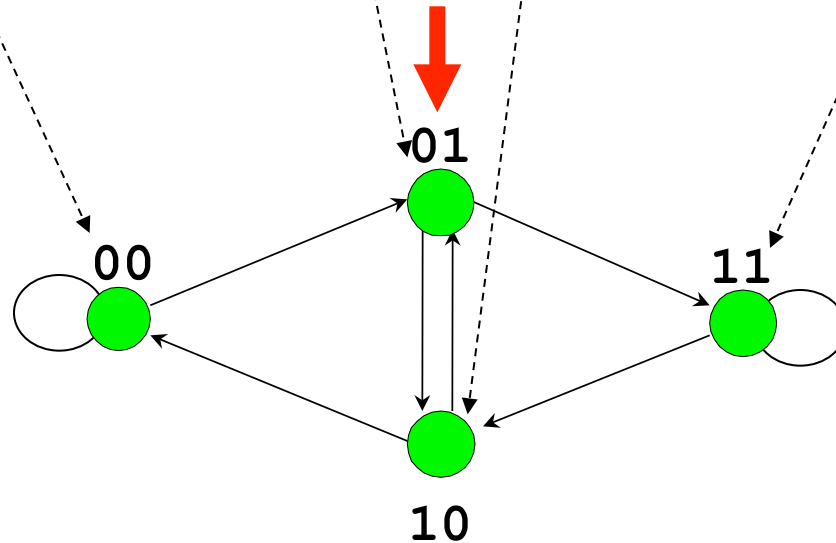
- Replace every edge of the  $N/2$ -node graph with a node
  - Edge  $(u_1u_2\dots u_r, u_2u_3\dots u_{r+1}) \rightarrow$  Node  $(u_1u_2\dots u_{r+1})$
- Insert a directed edge between pairs of nodes that correspond to consecutive directed edges in the  $N/2$ -node graph
  - Edge  $(u_1u_2\dots u_r, u_2u_3\dots u_{r+1})$  and Edge  $(u_2u_3\dots u_{r+1}, u_3u_4\dots u_{r+2})$
  - Replaced by  $(u_1u_2\dots u_{r+1}, u_2u_3\dots u_{r+2})$

# de Bruijn Networks: Construction Example

$N = 2$

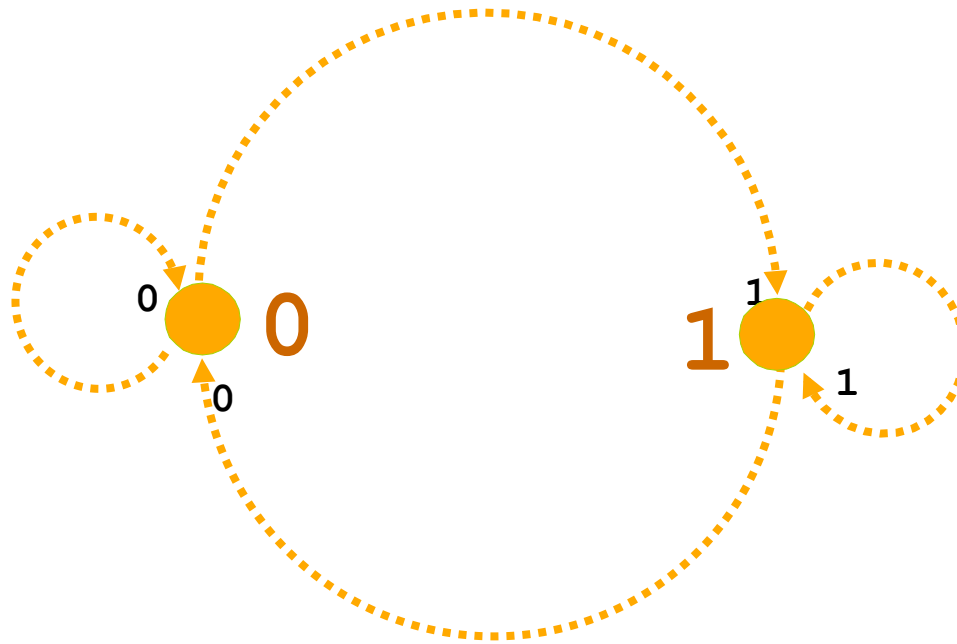


$N = 4$



# de Bruijn Networks: Construction Example

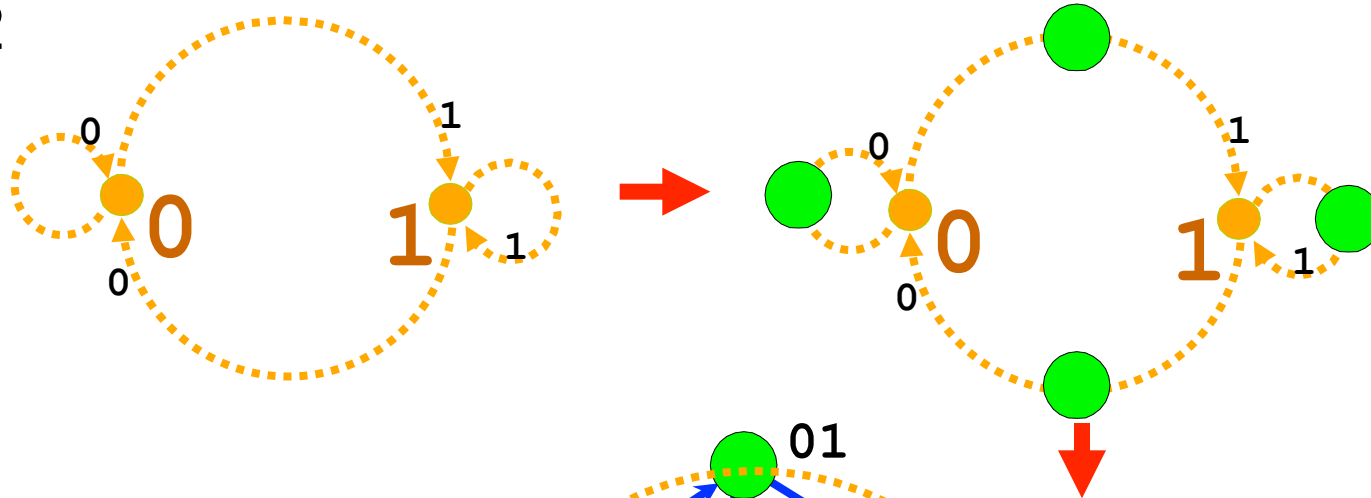
$N = 2$



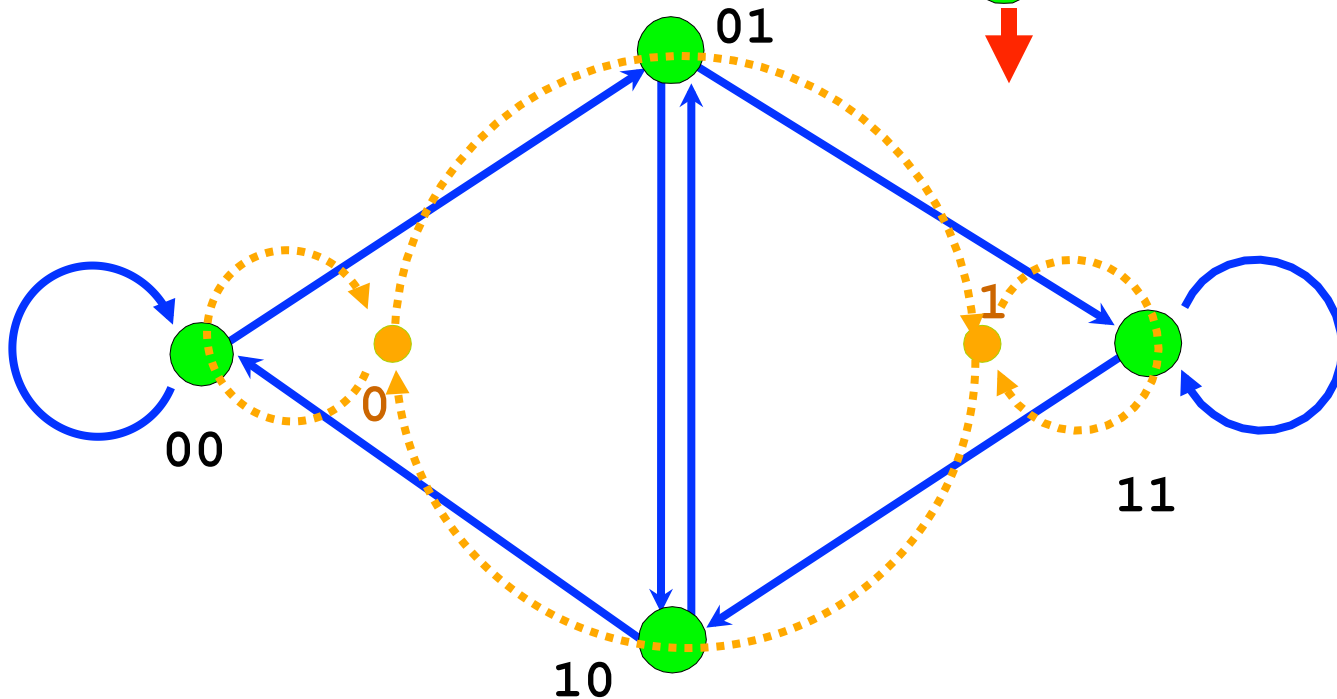


# de Bruijn Networks: Construction Example

$N = 2$



$N = 4$



## Scenarios:

- Controlled environments (provider defined topologies)

## Benefits

- Fixed and low node degree
- Fast routing:  $O(\log N)$

## Drawbacks

- Exponentially expandable
  - Network size is defined only for binary graphs
    - Node count: 1, 2, 4, 8, 16 ...  $2^i$
    - Not incrementally expandable
    - i.e. basic network cannot be defined for any arbitrary integer
- Cannot handle churn