

INFO8002

Large-Scale Data Systems

Exercise Session #4

Academic year 2021-2022



Reminder

REMINDER :

Distributed Hash Tables

- A **distributed hash table** (DHT) is a class of (decentralized?) distributed systems that provide a lookup services similar to a hash table.
 - Every node in a distributed hash table is responsible for a set of keys and their associated values.
 - The **key** is a unique identifier for its associated data value, obtained using a hashing function.
 - The **data values** can be any form of data.

REMINDER :

Distributed Hash Tables

Chord Algorithm: a protocol and algorithm for a peer-to-peer distributed hash table.

REMINDER :

Distributed Hash Tables

Chord Algorithm: a protocol and algorithm for finding a peer c such that $h(k) \in [id(c), id(c+1))$. **Not Really!**

- **Interface:**

- Support a single operation: **lookup(k)** → Return the ip of the host which hold the data associated to **k**.

- **Properties:**

- If **k** is stored on the DHT, a process will **eventually** find a node which stores **k**.
- **Termination**

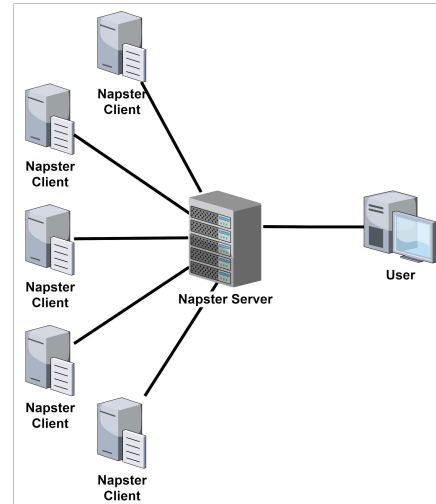


Help to locate where a resource is!

PROBLEM 1

Content Sharing System

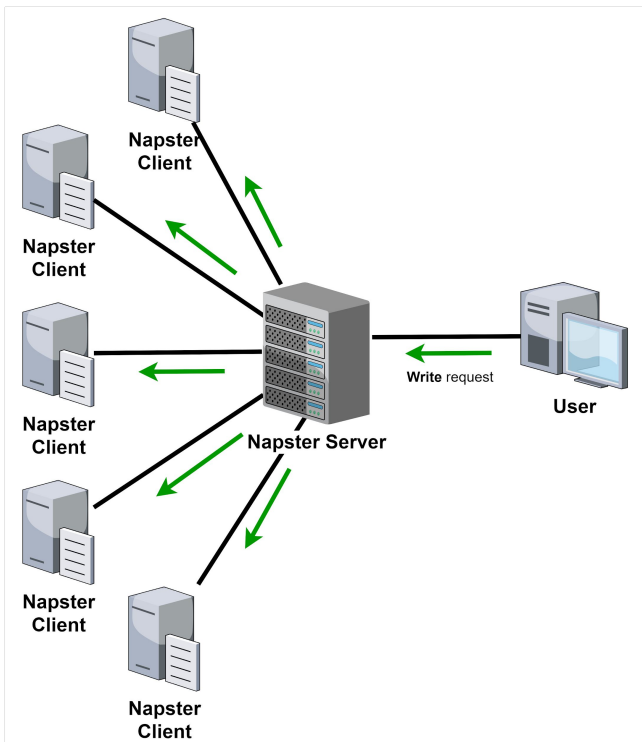
You are responsible for designing a system allowing the **storage** and **distribution** of content.



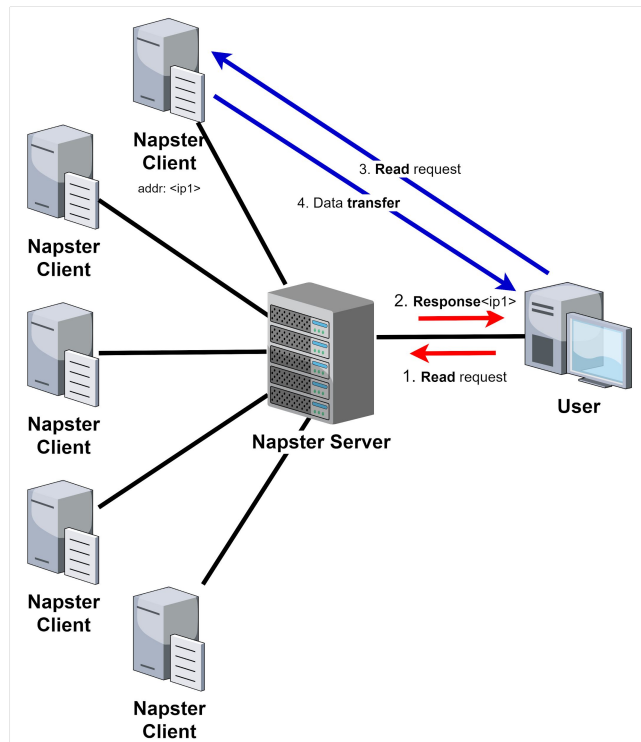
Specify an architecture for this distributed storage system and **provide** a pseudo-implementation using distributed hash tables.

Content Sharing System - 1st Design

Write query



Read query



Content Sharing System - 1st Design

Module Specification:

Module 1: Interface and properties of distributed storage

Module:

Name: *NapsterClientServerRegister*, **instance** *npr*.

Events:

Request: $\langle npr, \text{Read} \mid r, m \rangle$: Invokes a **read** operation on **m** consecutive registers starting on register **r**.

Request: $\langle npr, \text{Write} \mid v, r \rangle$: Invokes a **write** operation with value **v** starting on register **r**.

Indication: $\langle npr, \text{ReadReturn} \mid v \rangle$: Completes a **read** operation with return value **v**.

Indication: $\langle npr, \text{WriteReturn} \rangle$: Completes a **write** operation.

Properties:

NP1: Termination.

NP2: Validity.

Content Sharing System - 1st Design

Implementation

Algorithm 3:

Implements:

NapsterClientServerRegister, **instance** *npr*.

Uses:

(1, N)-NewRegularRegister, **instance** *onrr*.

upon event $\langle npr, \text{Init} \rangle$ **do**

pendingR := *pendingWr* := \emptyset ;

upon event $\langle npr, \text{Write} \mid v, r \rangle$ **do**

forall $v' \in v$ **do**

pendingWr := *pendingWr* $\cup \{r + \text{index}(v)\}$;

forall $v' \in v$ **do**

trigger $\langle onrr, \text{Write} \mid v', r + \text{index}(v) \rangle$;

upon event $\langle npr, \text{Read} \mid r, m \rangle$ **do**

ReadRet := $[0]^m$; *offset* := *r*;

for i **in** *range*(*m*) **do** *pendingR* := *pendingR* $\cup \{r + i\}$;

for i **in** *range*(*m*) **do** **trigger** $\langle onrr, \text{Read} \mid r + i \rangle$;

upon event $\langle onrr, \text{ReadReturn} \mid r, v \rangle$ **do**

pendingR := *pendingR* $\setminus \{r\}$;

ReadRet[*r*-*offset*] := *v*;

if *pendingR* $\subseteq \emptyset$ **then**

trigger $\langle np, \text{ReadReturn} \mid \text{ReadRet} \rangle$;

upon event $\langle onrr, \text{WriteReturn} \mid r \rangle$ **do**

pendingWr := *pendingWr* $\setminus \{r\}$;

if *pendingWr* $\subseteq \emptyset$ **then**

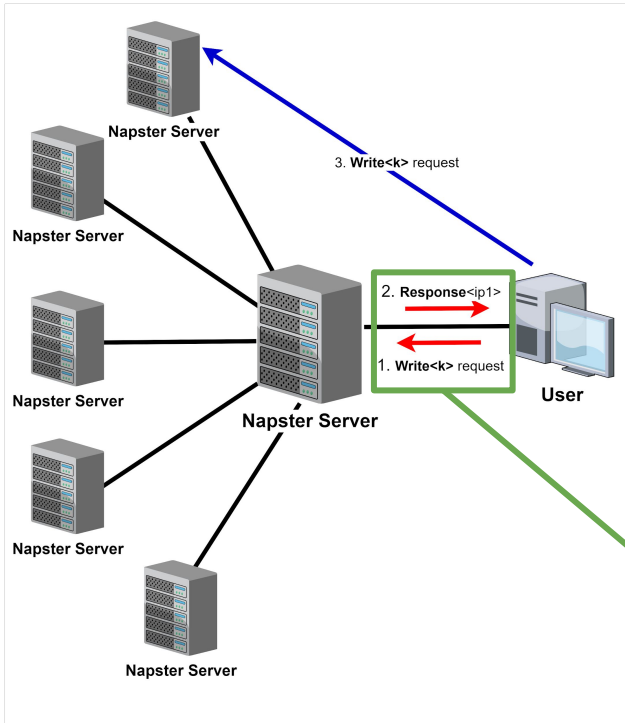
trigger $\langle onnrr, \text{Flush} \rangle$

upon event $\langle onnrr, \text{FlushReturn} \rangle$ **do**

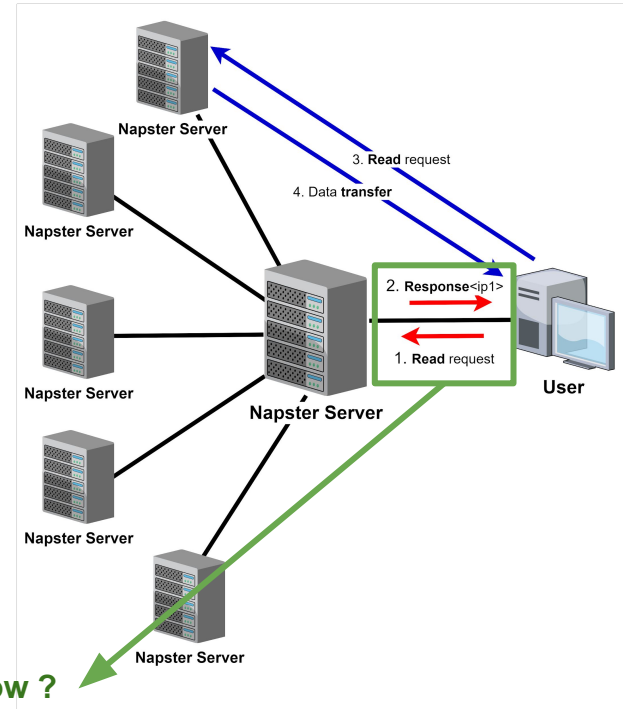
trigger $\langle np, \text{WriteReturn} \rangle$;

Content Sharing System - Updated

Write query



Read query



How ?

Content Sharing System - Updated

Module Specification:

Module 1: Interface and properties of distributed storage

Module:

Name: *NapsterServer*, **instance** *np*.

Events:

Request: $\langle np, \text{Read} \mid k \rangle$: Invokes a **read** operation of value with key **k**.

Request: $\langle np, \text{Write} \mid k, v \rangle$: Invokes a **write** operation of value **v** with key **k**.

Indication: $\langle np, \text{ReadReturn} \mid v \rangle$: Completes a **read** operation with return value **v**.

Indication: $\langle np, \text{WriteReturn} \rangle$: Completes a **write** operation.

Content Sharing System - Updated

Problems:

1. How can we ensure that operations terminates ?

Property n°1: "Termination"

"If a correct process invokes an operation, then the operation eventually completes."

2. How can we ensure that users receives a coherent response to their read request ?

Property n°2: "Validity"

"A get that is not concurrent with a put returns the last value written; a get that is concurrent with a put returns the last value written or the value concurrently written or no value."

Content Sharing System - Updated

Module Specification:

Module 1: Interface and properties of distributed storage

Module:

Name: *NapsterServer*, **instance** *np*.

Events:

Request: $\langle np, \text{Read} \mid k \rangle$: Invokes a **read** operation of value with key **k**.

Request: $\langle np, \text{Write} \mid k, v \rangle$: Invokes a **write** operation of value **v** with key **k**.

Indication: $\langle np, \text{ReadReturn} \mid v \rangle$: Completes a **read** operation with return value **v**.

Indication: $\langle np, \text{WriteReturn} \rangle$: Completes a **write** operation.

Properties:

NP1: Termination.

NP2: Validity.

Content Sharing System - Updated

Implementation

Algorithm 1:

Implements:

NapsterServer, **instance** *np*.

Uses:

Chord, **instance** *chord*.

NapsterClientServerRegister, **instance** *npr*.

upon event $\langle np, \text{Init} \rangle$ **do**

???

upon event $\langle np, \text{Write} \mid k, v \rangle$ **do**

???

upon event $\langle np, \text{Read} \mid k \rangle$ **do**

???

upon event $\langle chord, \text{lookupComplete} \mid k, ip, r, m \rangle$ **do**

???

upon event $\langle npr, \text{ReadReturn} \mid v \rangle$ **do**

???

upon event $\langle npr, \text{WriteReturn} \rangle$ **do**

???

Content Sharing System - Updated

Implementation

Algorithm 1:

Implements:

NapsterServer, instance *np*.

Uses:

Chord, instance *chord*.

NapsterClientServerRegister, instance *npr*.

upon event $\langle np, \text{Init} \rangle$ do

$pendingR := \emptyset;$
 $pendingWr := \emptyset;$

upon event $\langle np, \text{Write} \mid k, v \rangle$ do

$pendingWr := pendingWr \cup \{k : v\};$
trigger $\langle chord, \text{lookup} \mid k \rangle;$

upon event $\langle np, \text{Read} \mid k \rangle$ do

$pendingR := pendingR \cup \{k\};$
trigger $\langle chord, \text{lookup} \mid k \rangle;$

upon event $\langle chord, \text{lookupComplete} \mid k, ip, r, m \rangle$ do

if $k \in pendingR$ then

trigger $\langle npr, ip, \text{Read} \mid r, m \rangle;$
 $pendingR := pendingR \setminus \{k\};$

if $k \in pendingWr$ then

trigger $\langle npr, ip, \text{Write} \mid pendingWr[k], r \rangle;$
 $pendingWr := pendingWr \setminus \{k\};$

upon event $\langle npr, \text{ReadReturn} \mid v \rangle$ do

trigger $\langle np, \text{ReadReturn} \mid v \rangle;$

upon event $\langle npr, \text{WriteReturn} \rangle$ do

trigger $\langle np, \text{WriteReturn} \rangle;$

Content Sharing System - Updated

Chord vs Kademlia:

- **Kademlia** specifies how values should be stored & retrieved.
 - Resilient to node failures with persistent data storage.
 - Handling of nodes leaving/failing → straightforward (**by design**).
 - Provable properties in 1st paper → **Chord** needed further specifications¹.

- **Performance** comparisons of both algorithms are debatable depending on situations².

- **OTHER DESIGNS:**
 - CAN, Pastry, Tapestry → $O(\log(n))$ lookup (2001),
 - Koorde (2003) → $O(\log(n)/\log(\log(n)))$ lookup but **complex** to implement!

1. Zave, Pamela. "Using lightweight modeling to understand Chord." *ACM SIGCOMM Computer Communication Review* 42.2 (2012): 49-57.

2. Li, Jinyang, et al. "Comparing the performance of distributed hash tables under churn." *International Workshop on Peer-to-Peer Systems*. Springer, Berlin, Heidelberg, 2004.