# I 4-736 DISTRIBUTED SYSTEMS

RPC (KESDEN, SPRING 2019)

MANY THANKS TO PROFS. SRINI AND AGARWAL

... MOST OF THE PC SLIDES ARE FROM THEIR FALL 2016 15-440/640 LECTURE 6 SLIDE DECK

# A QUICK STORY: TWO ARMIES (AND ONE MESSENGER)



#### **COMMON COMMUNICATION PATTERN**



#### WRITING IT BY HAND ...

• E.g., if you had to write a, say, password cracker

Then wait for response, etc.

# <sup>6</sup> TODAY'S LECTURE

• RPC overview

• RPC challenges

• RPC other stuff

#### RPC – REMOTE PROCEDURE CALL: REPLACE COMMUNICATION VIA THE STACK WITH THE NETWORK

- A type of client/server communication
- Attempts to make remote procedure calls look like local ones



#### **RPC GOALS**

- Ease of programming
- Hide complexity
- Automates task of implementing distributed computation
- Familiar model for programmers (just make a function call)

Historical note: Seems obvious in retrospect, but RPC was only invented in the '80s. See

- Birrell & Nelson, "Implementing Remote Procedure Call" ... or
- Bruce Nelson, Ph.D. Thesis, Carnegie Mellon University: Remote Procedure Call., 1981 :)

# <sup>®</sup> REMOTE PROCEDURE CALL

- A remote procedure call makes a call to a remote service look like a local call
  - RPC makes transparent whether server is local or remote
  - RPC allows applications to become distributed transparently
  - RPC makes architecture of remote machine transparent

## BUT IT'S NOT ALWAYS SIMPLE

- Calling and called procedures run on different machines, with different address spaces
  - And perhaps different environments .. or operating systems ..
- Must convert to local representation of data
- Machines and network can fail

# <sup>11</sup> STUBS: OBTAINING TRANSPARENCY

- Compiler generates from API stubs for a procedure on the client and server
- Client stub
  - Marshals arguments into machine-independent format
  - Sends request to server
  - Waits for response
  - <u>Unmarshals</u> result and returns to caller
- Server stub
  - <u>Unmarshals</u> arguments and builds stack frame
  - Calls procedure
  - Server stub marshals results and sends reply

## WRITING IT BY HAND - (AGAIN...)

• E.g., if you had to write a, say, password cracker

Then wait for response, etc.

## MARSHALING AND UNMARSHALING

- (From example) htonl() -- "host to network-byte-order, long".
  - network-byte-order (big-endian) standardized to deal with cross-platform variance
- Note how we arbitrarily decided to send the string by sending its length followed by L bytes of the string? That's marshalling, too.
- Floating point...
- Nested structures? (Design question for the RPC system do you support them?)
- Complex datastructures? (Some RPC systems let you send lists and maps as first-order objects)



- a) Original message on x86 (Little Endian)
- b) The message after receipt on the SPARC (Big Endian)
- c) The message after being inverted. The little numbers in boxes indicate the address of each byte

## "STUBS" AND IDLS

- RPC stubs do the work of marshaling and unmarshaling data
- But how do they know how to do it?
- Typically: Write a description of the function signature using an IDL -- interface definition language.
  - Lots of these. Some look like C, some look like XML, ... details don't matter much.

# <sup>16</sup> PASSING VALUE PARAMETERS (I)



3. Message is sent across the network

• The steps involved in a doing a

remote computation through RPC.

# <sup>18</sup> REMOTE PROCEDURE CALLS (2)

• A remote procedure call occurs in the following steps (continued):

6.The server does the work and returns the result to the stub.
7.The server stub packs it in a message and calls its local OS.
8.The server's OS sends the message to the client's OS.
9.The client's OS gives the message to the client stub.
10.The stub unpacks the result and returns to the client.

# <sup>19</sup> PASSING REFERENCE PARAMETERS

- Replace with pass by copy/restore
- Need to know size of data to copy
  - Difficult in some programming languages
- Solves the problem only partially
  - What about data structures containing pointers?
  - Access to memory in general?

#### <sup>20</sup> TWO STYLES OF RPC IMPLEMENTATION

- Shallow integration. Must use lots of library calls to set things up:
  - How to format data
  - Registering which functions are available and how they are invoked.
  - Often seen in heterogeneous uses, i.e., cross-environment, Java to C
- Deep integration.
  - Data formatting done based on type declarations
  - (Almost) all public methods of object are registered.
  - Often seen in homogeneous uses, e.g. Go or Java.

# <sup>21</sup> TODAY'S LECTURE

• RPC overview

• RPC challenges

• RPC other stuff

# <sup>22</sup> RPC VS. LPC

- 3 properties of distributed computing that make achieving transparency difficult:
  - Partial failures
  - Latency
  - Memory access

#### **RPC FAILURES**

- Request from cli  $\rightarrow$  srv lost
- Reply from srv  $\rightarrow$  cli lost
- Server crashes after receiving request
- Client crashes after sending request

# <sup>24</sup> PARTIAL FAILURES

- In local computing:
  - if machine fails, application fails
- In distributed computing:
  - if a machine fails, part of application fails
  - one cannot tell the difference between a machine failure and network failure
- How to make partial failures transparent to client?

# <sup>25</sup> STRAWMAN SOLUTION

- Make remote behavior identical to local behavior:
  - Every partial failure results in complete failure
    - You abort and reboot the whole system
  - You wait patiently until system is repaired
- Problems with this solution:
  - Many catastrophic failures
  - Clients block for long periods
    - System might not be able to recover

# <sup>26</sup> REAL SOLUTION: BREAK TRANSPARENCY

- Possible semantics for RPC:
  - Exactly-once
    - Impossible in practice
  - At least once:
    - Only for idempotent operations
  - At most once
    - Zero, don't know, or once
  - Zero or once
    - Transactional semantics

#### **EXACTLY-ONCE?**

- Sorry no can do in general.
- Imagine that message triggers an external physical thing (say, a robot fires a nerf dart at the professor)
- The robot could crash immediately before or after firing and lose its state. Don't know which one happened. Can, however, make this window very small.

## **REAL SOLUTION: BREAK TRANSPARENCY**

- **<u>At-least-once</u>**: Just keep retrying on client side until you get a response.
  - Server just processes requests as normal, doesn't remember anything. Simple!
- **<u>At-most-once</u>**: Server might get same request twice...
  - Must re-send previous reply and not process request (implies: keep cache of handled requests/responses)
  - Must be able to identify requests
  - Strawman: remember all RPC IDs handled.
    - $\rightarrow$  Ugh! Requires infinite memory.
  - Real: Keep sliding window of valid RPC IDs, have client number them sequentially.

## **IMPLEMENTATION CONCERNS**

- As a general library, performance is often a big concern for RPC systems
- Major source of overhead: copies and marshaling/unmarshaling overhead
- Zero-copy tricks:
  - Representation: Send on the wire in native format and indicate that format with a bit/byte beforehand. What does this do? Think about sending uint32 between two little-endian machines
  - Scatter-gather writes (writev() and friends)

# DEALING WITH ENVIRONMENTAL DIFFERENCES

- If my function does: read(foo, ...)
- Can I make it look like it was really a local procedure call??
- Maybe!
  - Distributed filesystem...
- But what about address space?
  - This is called distributed shared memory
  - People have kind of given up on it it turns out often better to admit that you're doing things remotely

# <sup>31</sup> TODAY'S LECTURE

• RPC overview

• RPC challenges

• RPC other stuff

# <sup>32</sup> ASYNCHRONOUS RPC (I)



• The interaction between client and server in a traditional RPC.

# <sup>33</sup> ASYNCHRONOUS RPC (2)



• The interaction using asynchronous RPC.

# <sup>34</sup> ASYNCHRONOUS RPC (3)



• A client and server interacting through

two asynchronous RPCs.

## <sup>35</sup> USING RPC

- How about a distributed bitcoin miner.
- Three classes of agents:
  - I. Request client. Submits cracking request to server. Waits until server responds.
  - 2. Worker. Initially a client. Sends join request to server. Now it should reverse role & become a server. Then it can receive requests from main server to attempt cracking over limited range.
  - 3. Server. Orchestrates whole thing. Maintains collection of workers. When receive request from client, split into smaller jobs over limited ranges. Farm these out to workers. When finds bitcoin, or exhausts complete range, respond to request client.

## 36 USING RPC

- Request→Server→Response:
  - Classic synchronous RPC
- Worker→Server.
  - Synch RPC, but no return value.
  - "Í'm a worker and I'm listening for you on host XXX, port YYY."
- Server→Worker.
  - Synch RPC?
  - No that would be a bad idea. Better be Asynch.
  - Otherwise, it would have to block while worker does its work, which misses the whole point of having many workers.

# <sup>37</sup> BINDING A CLIENT TO A SERVER

- Registration of a server makes it possible for a client to locate the server and bind to it
- Server location is done in two steps:
  - Locate the server's machine.
  - Locate the server on that machine.



# <sup>38</sup> SUMMARY: EXPOSE REMOTENESS TO CLIENT

• Expose RPC properties to client, since you cannot hide them

- Application writers have to decide how to deal with partial failures
  - Consider: E-commerce application vs. game

## <sup>39</sup> IMPORTANT LESSONS

- Procedure calls
  - Simple way to pass control and data
  - Elegant transparent way to distribute application
  - Not only way...
- Hard to provide true transparency
  - Failures
  - Performance
  - Memory access
  - Etc.
- How to deal with hard problem  $\rightarrow$  give up and let programmer deal with it
  - "Worse is better"