

# JGroups

COMUNICAZIONE DI GRUPPO

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Panoramica

JGroups è un framework che offre servizi di comunicazione di gruppo (jgroups.org) fondata dell'application server JBoss (jboss.org)

Servizi di comunicazione:

- Consegna garantita
- Gestione dei frammenti
- Ordinamento
- Consegna atomica

	Inaffidabile	Affidabile
unicast	UDP	TCP
multicast	IP-Multicast	JGroups

Servizi di gruppo:

- Stato condiviso
- Gestione membri
- Notifica di ingresso, uscita, fallimento dei membri

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Lo stack JGroups

JGroups ha un'architettura a stack

È possibile aggiungere strati (protocolli) secondo necessità

Molti sono già presenti nella libreria

- Diversi protocolli di trasporto (UDP IP-Multicast, TCP, ...)
- Affidabilità
- Ordinamento
- Compressione
- Crittografia
- Gestione del flusso (per bilanciare consumatori lenti e produttori veloci)
- Gestione dei fallimenti/partizionamenti

Ma è possibile utilizzarne di terze parti (o propri)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

# Architettura di gruppo

Membri di un gruppo (cluster) sono pari (ma con ruoli differenti)

Eleggono un coordinatore (solitamente il più anziano)

Il coordinatore

- Tiene traccia dei membri del gruppo
- Mantiene lo stato condiviso (e gestisce la replicazione)
- Gestisce la comunicazione di gruppo
  - Ordinamento totale (un nodo chiede al coordinatore di distribuire un messaggio)
  - Garanzia di consegna (ack/retry...)

Fallback secondo «anzianità» nel gruppo

(almeno nella configurazione «base»...vedremo come rendere le cose più interessanti)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

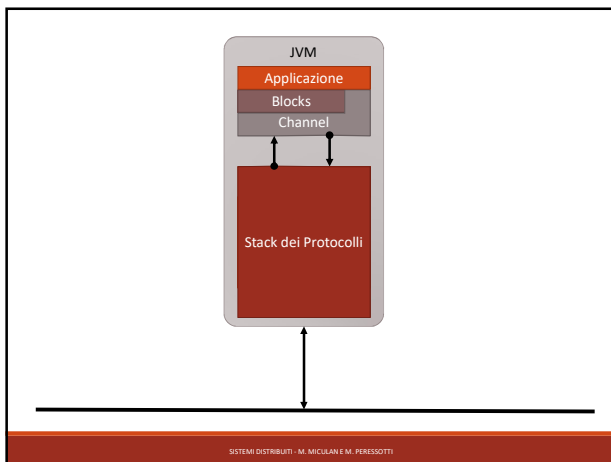
---

---

---

---

---



SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

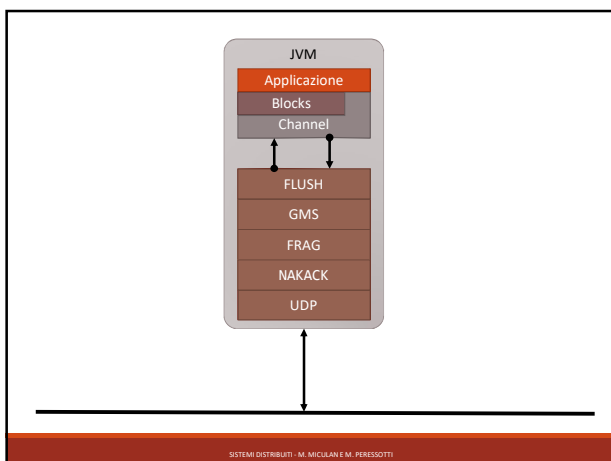
---

---

---

---

---



SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

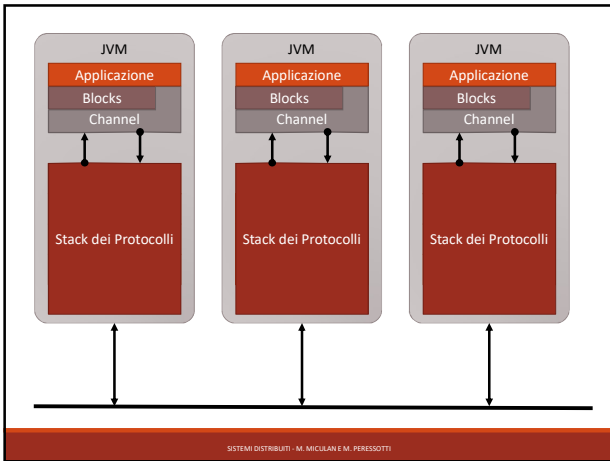
---

---

---

---

---




---

---

---

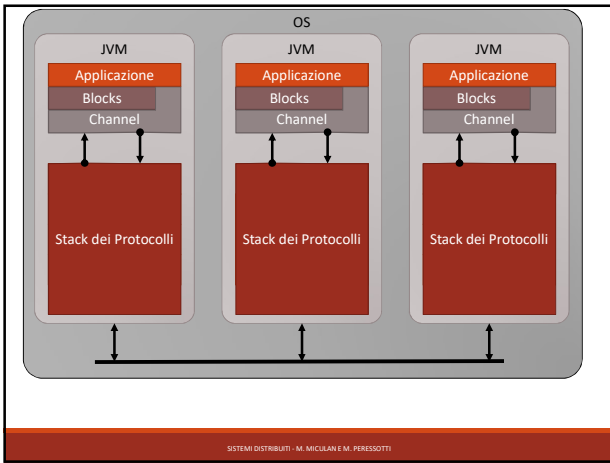
---

---

---

---

---




---

---

---

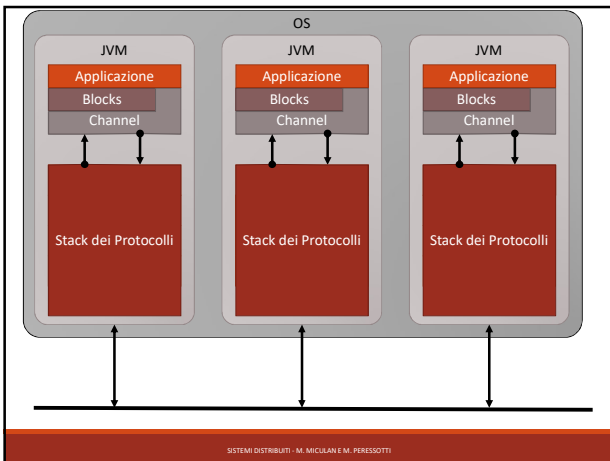
---

---

---

---

---




---

---

---

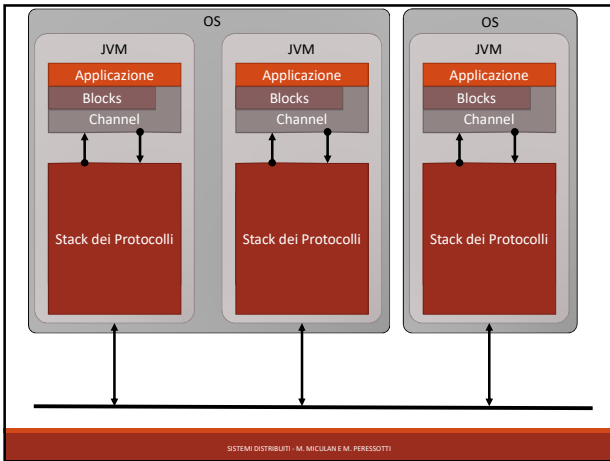
---

---

---

---

---




---

---

---

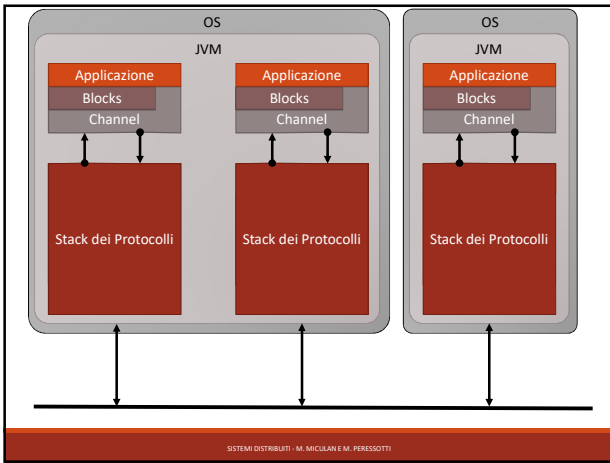
---

---

---

---

---




---

---

---

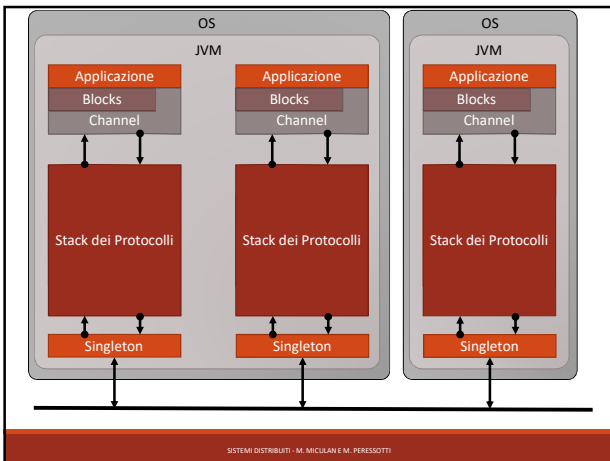
---

---

---

---

---




---

---

---

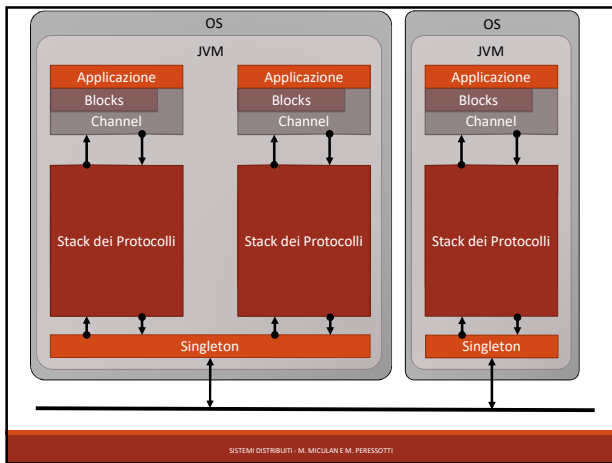
---

---

---

---

---




---

---

---

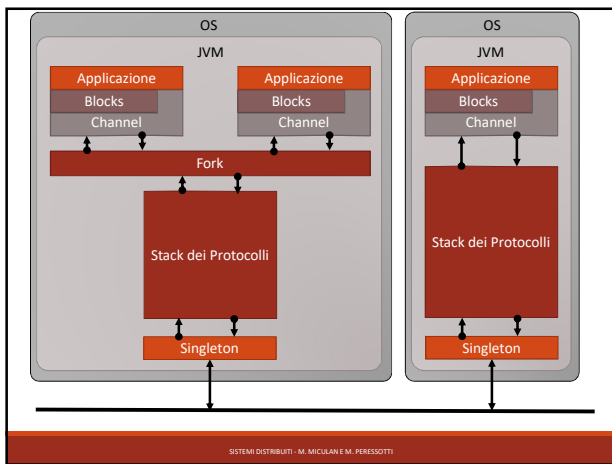
---

---

---

---

---




---

---

---

---

---

---

---

---

## JChannel

La classe JChannel costituisce l'API principale della libreria

- Connessione a un gruppo (cluster)
- Invio/ricezione dei messaggi
- Eventi ed ascoltatori (e.g. ingresso uscita di nodi dal gruppo)

La classe Message definisce i messaggi

- Il payload deve essere serializzabile
- Manca polimorfismo parametrico (generics)

I messaggi possono essere inviati e ricevuti solamente da membri dello stesso gruppo

I membri del gruppo sono noti e reperibili mediante viste (View)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## JChannel

Creazione e connessione ad un canale:

```
channel=new JChannel(); //configurazione predefinita udp.xml
channel.connect("NomeGruppo");
```

Invio di un messaggio:

```
channel.send(new Message(
    null, // destinatario, null -> gruppo
    null, // indirizzo risposte, null -> "questa" istanza
    "ciao" // payload
));
```

Disconnessione e chiusura

```
channel.disconnect();
channel.close();
```

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

---

---

## Ricezione

La ricezione di messaggi ed eventi è demandata a implementazioni di **Receiver** (e.g. ReceiverAdapter)

```
public class Chat extends ReceiverAdapter {
    ...
    @Override
    public void receive(Message msg){
        Object payload = msg.getObject();
        ...
    }
}
```

Infine, il ricevitore va registrato presso il canale:

```
channel.connect("NomeGruppo");
...
channel.setReceiver(this);
```

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

---

---

## Viste e partizionamenti

La lista dei membri di un gruppo è reperibile mediante

```
View view = channel.getView();
```

Le viste sono cristallizzate i.e. non sono aggiornate automaticamente. Gli aggiornamenti sono notificati al receiver mediante la callback:

```
viewAccepted(View newView)
```

I gruppi possono essere partizionati a seguito di fallimenti della rete sottostante (non è possibile distinguere fallimenti di link e nodi)

Le parti evolvono indipendentemente rendendo necessario il merge (se e quando il link sarà ripristinato)

Il merge è demandato all'utente cui è fornita una **MergeView** (i.e. un albero di view)

(Esempio...)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

---

---

## Stato «condiviso»

I Receiver possono esportare/importare uno stato su stream

- getState
- setState

I membri del gruppo possono richiedere lo stato degli altri membri:

```
channel.getState(member, timeout);
```

Non possono modificarlo senza collaborazione

Lo stato **non è condiviso**, ogni istanza ha la propria «visione»

Uso tipico:

- stato replicato,
- aggiornato alla ricezione dei messaggi
- i nuovi membri copiano lo stato dal coordinatore

(esempio...)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Ascoltatori

JChannel consente di registrare ascoltatori per vari eventi:

- `MessageListener`
  - void `receive(Message msg)`
  - void `getState(OutputStream output)` throws Exception
  - void `setState(InputStream input)` throws Exception
- `MembershipListener`
  - void `viewAccepted(View new_view)` cambiamenti nella vista
  - void `suspect(Object suspected_member)` fallimenti sospettati
  - void `block()` notifica la sospensione dell'invio di messaggi (e.g. FLUSH e STATE\_TRANSFER) breve!
  - void `unblock()` ripristino invio di messaggi

`ReceiverAdapter` implementa `Receiver` estende questi ascoltatori

**Attenzione** alle dipendenze tra questi eventi o si rischia di bloccare lo stack (e.g. invio/ricezione di messaggi durante `viewAccepted` in presenza di FLUSH)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Indirizzi

Gli indirizzi JGroups sono UUIDs

Nasconde gli indirizzi dello stack sottostante

È possibile personalizzare gli indirizzi:

1. Implementare `org.jgroups.stack.AddressGenerator`
2. Registrazione `Jchannel.setAddressGenerator(AddressGenerator)`

Alcuni accorgimenti:

- Gli indirizzi devono implementare `Address`
- Gli indirizzi devono estendere `UUID`
- Non sovrascrivere `compare()`, `equals()` e `hashCode`
- Registrazione (della classe, non dell'istanza)
  - `ClassConfigurator.getInstance().put(1234, AwesomeAddress.class)`
  - `ig-magic-map.xml`

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Indirizzi

Perché provvedere indirizzi diversi?

- (perché possiamo)
- UUID sono unici ma non riflettono e.g. strutture amministrative
- UUID sono poco «leggibili» (mnemonici, significativi, ...)

Perché provvedere generatori diversi?

- Aggiunta di informazioni
- Controllo sull'assegnazione degli indirizzi

Alternative (o basi per *TheChefsChoiceOfTheDayAddress*) in JGroups:

- AdditionalDataUUID: campo binario
- PayloadUUID: campo testuale
- TopologyUUID: campi sito, rack, macchina

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

---

---

## Stack dei protocolli

JGroups ha uno stack di protocolli altamente configurabile.

La classe `ProtocolStack` consente di manipolare lo stack:

- `get/find/add/insert/remove/replaceProtocol`
- `printProtocolSpec`
- `dumpStats`
- `down/up(Event)`

Il metodo `channel.getProtocolStack()` restituisce lo stack del canale.

(esempi...)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

---

---

## Configurare lo stack: java

```
JChannel channel=new JChannel(false);
ProtocolStack stack=new ProtocolStack();
channel.setProtocolStack(stack);
stack.addProtocol(new UDP());
    .addProtocol(new PING());
    .addProtocol(new MERGE2());
    .addProtocol(new FD_SOCK());
    .addProtocol(new FD_ALL()
        .setValue("timeout", 12000)
        .setValue("interval", 3000));
    .addProtocol(new VERIFY_SUSPECT());
    .addProtocol(new BARRIER());
    .addProtocol(new NAKACK());
    .addProtocol(new UNICAST2());
    .addProtocol(new STABLE());
    .addProtocol(new GMS());
    .addProtocol(new UFC());
    .addProtocol(new MFC());
    .addProtocol(new FRAG2());
stack.init();
```

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

---

---

## Configurare lo stack: xml

```
<config xmlns="urn:jgroups"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:schemaLocation="urn:jgroups
http://www.jgroups.org/schema/JGroups-3.4.xsd">
  <UDP ... />
  <PING ... />
  <MERGE2 ... />
  <FD_SOCKET/>
  <FD_ALL "timeout"=12000 "interval"=3000/>
  <VERIFY_SUSPECT ... />
  <BARRIER />
  <pbcast.NAKACK ... />
  <UNICAST2 ... />
  <pbcast.STABLE .../>
  <pbcast.GMS .../>
  <UFC .../>
  <MFC .../>
  <FRAG2 ... />
</config>
```

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Configurazioni ed estensioni

Molte configurazioni sono distribuite assieme al pacchetto JGroups  
• E.g. udp.xml, tcp.xml

Molti (ma proprio molti) protocolli e servizi in org.jgroups.protocols  
• ([www.jgroups.org/manual/html/protlist.html](http://www.jgroups.org/manual/html/protlist.html))

Se questi non dovessero bastare, le terze parti sono supportate

Scrivere un protocollo:

1. Estendere org.jgroups.protocols.Protocol (o sottoclassi)
2. Fornire un nome univoco
3. Implementare i metodi up() e down()
4. (Eventuali header devono avere un nome univoco)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

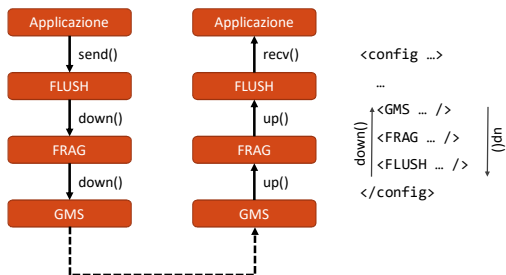
---

---

---

---

## Up(side)-down



SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Esempio: THROTTLER

```
public class THROTTLER extends Protocol{
    public Object down(final Event evt) {
        if (isMessage(evt) & coinFlip())
            return null;
        else
            return down_prot.down(evt);
    }

    public Object up(final Event evt) {
        if (isMessage(evt) & coinFlip())
            return null;
        else
            return up_prot.up(evt);
    }
}
```

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Header

I messaggi possono essere decorati da header

Veicolano informazione usata dai protocolli

È buona educazione **non toccare** gli header altrui...

...ma definire i propri (salvo casi marginali)

1. Estendere **Header** (almeno un costruttore vuoto)
2. Implementare **Streamable** (marshaling)
3. Nome univoco (e.g. annidarlo nella classe del protocollo)
4. Implementare **size()** (byte prodotti durante il marshaling)
5. Registrazione (java o xml)
  - ClassConfigurator.getInstance().put(1234, AwesomeHeader.class)
  - jg-magic-map.xml

Message espone putHeader(), getHeader() removeHeader()

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Giusto per farsi un'idea

Transport <ul style="list-style-type: none"><li>◦ UDP, TCP, TUNNEL</li></ul>	Group membership <ul style="list-style-type: none"><li>◦ GMS</li></ul>
Membership discovery <ul style="list-style-type: none"><li>◦ PING, TCPING, TCPGOSSIP, MPING, FILE_PING, JDBC_PING, BPING, RACKSPACE_PING, S3_PING, GOOGLE_PING, AWS_PING, PDC</li></ul>	Flow control <ul style="list-style-type: none"><li>◦ FC, MFC, UFC</li></ul>
Partition recovery <ul style="list-style-type: none"><li>◦ MERGE, MERGE2, MERGE3</li></ul>	Fragmentation <ul style="list-style-type: none"><li>◦ FRAG, FRAG2</li></ul>
Failure detection <ul style="list-style-type: none"><li>◦ FD, FD_ALL, FD_ALL2, FD_SOCK, FD_PING, FD_HOST, VERIFY_SUSPECT</li></ul>	Group ordering <ul style="list-style-type: none"><li>◦ SEQUENCER, TOA</li></ul>
Reliable message transmission <ul style="list-style-type: none"><li>◦ NAKACK, NAKACK2, UNICAST, UNICAST2, UNICAST3, RSVP</li></ul>	State transfer <ul style="list-style-type: none"><li>◦ STATE_TRANSFER, STATE, STATE_LOCK, BARRIER, FLUSH</li></ul>
Message stability <ul style="list-style-type: none"><li>◦ STABLE</li></ul>	Misc <ul style="list-style-type: none"><li>◦ STAT, ENCRYPT, AUTH, COMPRESS, SCOPE, RELAY, RELAY2, DAISYCHAIN, RATE_LIMITER, CENTRAL_LOCK, PEER_LOCK, COUNTER, CENTRAL_EXECUTOR, SUPERVISOR, FORK</li></ul>

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Esercizi (spunti)

1. Servizi replicati: client/server con server replicato
  1. Interazione client/server a scelta e.g. socket, RPC, RMI, REST etc.
  2. Diversificare i protocolli (server RMI, server REST etc.)
  3. Server «attivabili»
2. Cluster di computazione basato su task
  1. Task sono sottoposti al cluster
  2. I task sono distribuiti internamente ai membri del cluster
  3. Assunzione: i task terminano
  4. Ogni task sottoposto dovrà essere completato
3. Aggiungere priorità, tipologie, time-slot, dipendenze tra (sub)task
4. Gestione partizionamenti/merge in presenza di stato
5. Scope e bridge tra clusters (RELAY and SCOPE)
6. Giocare con lo stack (e.g. negoziare QoS)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## Blocchi

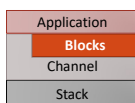
Formano uno strato aggiuntivo tra applicazione e canale

Offrono

- astrazione sui canali (che sono costruiti **socket-like** asincroni)
- Interfacce più articolate
- Servizi aggiuntivi (e.g. Stato, RPC)

Alcuni esempi

- Message dispatcher
- Rpc dispatcher
- Replicated hashmap
- Replicated cache
- Cluster-wide lock/counter/task/transaction



SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

## MessageDispatcher

Correlazione bloccante (e non) tra richiesta e risposta

Diverse tipologie di richiesta:

- GET\_ALL
- GET\_NONE
- GET\_FIRST
- GET\_MAJORITY

Esempio: richiesta di voto, servizio

Ulteriori opzioni:

- Timeout
- Unicast
- Filtri (e.g. invalida la richiesta di voto se c'è un veto)
- Exclusion list
- Scope (nel caso sia abilitato l'invio concorrente, cf protocollo SCOPE)
- Flags (gli stessi dei messaggi)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

# MessageDispatcher

Creazione di un dispatcher

- `MessageDispatcher disp = new MessageDispatcher(channel);`

Ascoltatore per richieste al dispatcher:

- `disp.setRequestHandler(new RequestHandler(){  
 public Object handle(Message msg) throws Exception {  
 ...  
 }  
});`

Invio di richieste:

- `castMessage(...)`
- `castMessageWithFuture(...)`
- `sendMessage(...)`
- `sendMessage(...)`

(Esempio: Poll)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---

# Esercizi (spunti)

7. Utilizzare i blocchi nei casi precedenti
  - (Non tutti consentono di lavorare con il canale sottostante, FORK?)
  - Voto nella chat
  - Risoluzione merge per votazione (dei processi)
  - Repliche con consistenza debole (letture e scritture per voto)
  - (La votazione è uno spunto, ci sono tanti altri problemi interessanti...)
  - Gestione dei lock (e.g. scrittura di risorse replicate)

SISTEMI DISTRIBUITI - M. MICULAN E M. PERESSOTTI

---

---

---

---

---

---

---

---