

PDV 12 2019/2020

# Algoritmus Raft

Michal Jakob

[michal.jakob@fel.cvut.cz](mailto:michal.jakob@fel.cvut.cz)

Centrum umělé inteligence, katedra počítačů, FEL ČVUT





# Algorithmus RAFT

# Model a požadavky

## Asynchronní systém se selháními

- procesy mohou **havarovat** (fail-stop, tj. nikoliv byzantsky)
- zprávy se mohou **ztrácet** (ale dodržují pořadí → nedokonalý FIFO kanál)

Algoritmus pro konsensus musí garantovat **bezpečnost** a měl by maximalizovat **živost** (dostupnost)

- obojí garantovat nelze ← **FLP teorém**

# Přístupu k problému konsensu

## **Symetrický/**bez lídra

- všechny servery mají stejnou roli
- klienti mohou kontaktovat kterýkoliv server

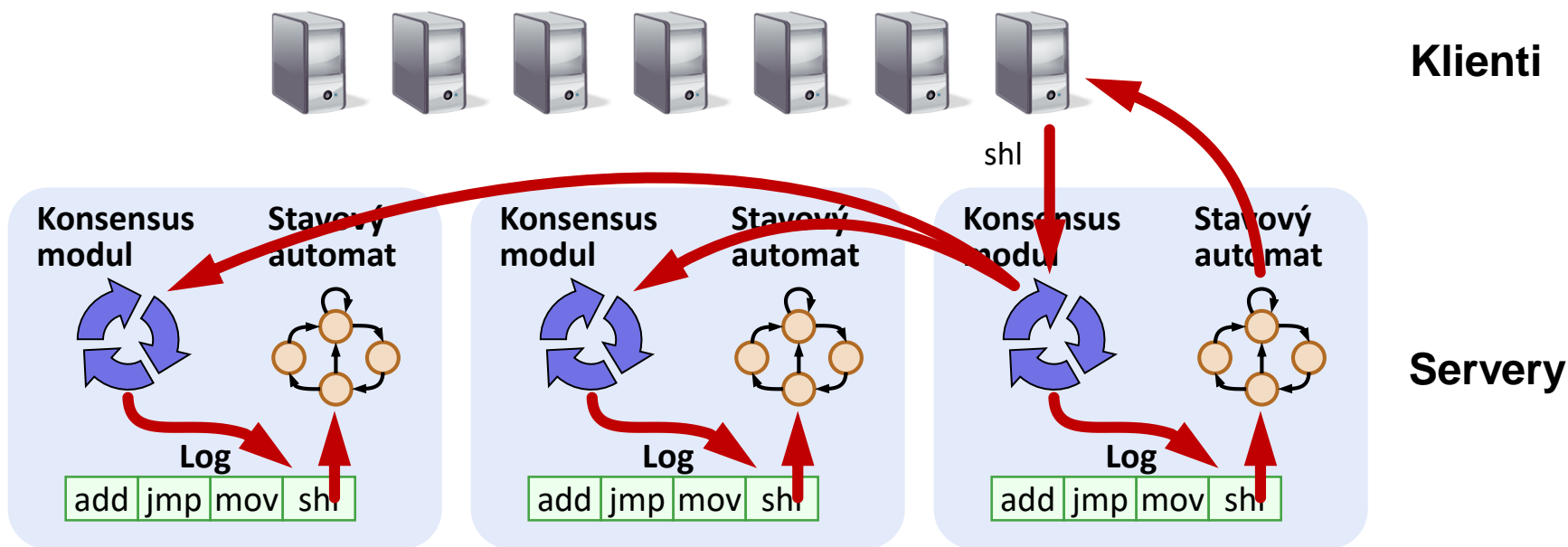
## **Asymetrický/**s lídrem

- v každém okamžiku je jeden server lídrem a ostatní přijímají jeho rozhodnutí
- klienti komunikují s lídrem

## Raft využívá lídra – výhody:

- **dekomponuje** problém na 1) **běžný chod** a 2) **změny lídra**
- **zjednodušuje** běžný chod (nedochází ke konfliktům)
- **efektivnější** než symetrické přístupy bez lídra (selhání lídra jsou v praxi vzácná)

# Cíl: Replikovaný log



Replikovaný log → **replikovaný stavový automat**

- Všechny procesy vykonávají příkazy ve stejném pořadí

Konsensus modul zajišťuje správnou **replikaci logu** a rozhoduje, kdy mohou být příkazy vykonány.

Zpracování požadavků klientů postupuje pokud je **nadpoloviční většina** serverů aktivních.

# Přehled Raftu

## 1. Volba lídra

- volba jednoho ze serveru jako lídra
- detekce selhání a vyvolání volby nového lídra

## 2. Běžný chod (základní replikace logu)

## 3. Bezpečnost a konzistence po změně lídra

## 4. Neutralizace starých lídrů

## 5. Interakce s klienty

## 6. (Rekonfigurace)



# Volba lídra

# Stavy serveru<sup>1</sup>

V každém okamžiku je každý server v právě **jednom stavu**:

## Lídr

obsluhuje  
požadavky klientů  
a replikuje log

## Následovník

pasivní – pouze  
reagují na zprávy  
od jiných serverů

## Kandidát

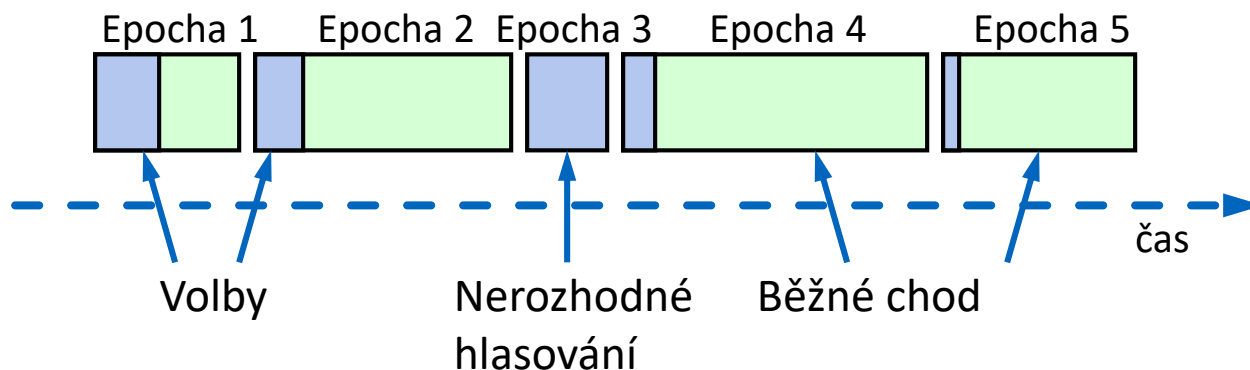
přechodná role  
v průběhu volby  
lídra

Běžný chod: 1 lídr ,  $N - 1$  následovníků

<sup>1</sup>Skupinu procesů, které se účastní konsensu, budeme označovat jako servery – pro lepší odlišení od procesů, které běží na klientských počítačích a které se konsensu neúčastní.



# Epochy (volební období)



Čas je rozdělen do **epoch** (~logický čas): každá epocha má své **číslo**, čísla jsou **inkrementována** a nikdy nejsou znovu použita. Každý server si (*persistentně!*) udržuje číslo **aktuální epochy**.

Epochy mohou mít dvě části

- **Volby** (buď nedopadnou nebo vyústí ve zvolení právě jednoho lídra)
- **Běžný chod** pod jedním zvoleným lídrem

Maximálně jeden lídr v každé epoše; některé epochy ale nemají lídra (neúspěšné volby).

Epochy slouží k identifikování **zastaralých informací**.

# Raft Protocol Summary

## Followers

- Respond to RPCs from candidates and leaders.
- Convert to candidate if election timeout elapses without either:
  - Receiving valid AppendEntries RPC, or
  - Granting vote to candidate

## Candidates

- Increment currentTerm, vote for self
- Reset election timeout
- Send RequestVote RPCs to all other servers, wait for either:
  - Votes received from majority of servers: become leader
  - AppendEntries RPC received from new leader: step down
  - Election timeout elapses without election resolution: increment term, start new election
- Discover higher term: step down

## Leaders

- Initialize nextIndex for each to last log index + 1
- Send initial empty AppendEntries RPCs (heartbeat) to each follower; repeat during idle periods to prevent election timeouts
- Accept commands from clients, append new entries to local log
- Whenever last log index  $\geq$  nextIndex for a follower, send AppendEntries RPC with log entries starting at nextIndex, update nextIndex if successful
- If AppendEntries fails because of log inconsistency, decrement nextIndex and retry
- Mark log entries committed if stored on a majority of servers and at least one entry from current term is stored on a majority of servers
- Step down if currentTerm changes

## Persistent State

Each server persists the following to stable storage synchronously before responding to RPCs:

<b>currentTerm</b>	latest term server has seen (initialized to 0 on first boot)
<b>votedFor</b>	candidateId that received vote in current term (or null if none)
<b>log[]</b>	log entries

## Log Entry

<b>term</b>	term when entry was received by leader
<b>index</b>	position of entry in the log
<b>command</b>	command for state machine

## RequestVote RPC

Invoked by candidates to gather votes.

### Arguments:

<b>candidateId</b>	candidate requesting vote
<b>term</b>	candidate's term
<b>lastLogIndex</b>	index of candidate's last log entry
<b>lastLogTerm</b>	term of candidate's last log entry

### Results:

<b>term</b>	currentTerm, for candidate to update itself
<b>voteGranted</b>	true means candidate received vote

### Implementation:

1. If term > currentTerm, currentTerm  $\leftarrow$  term (step down if leader or candidate)
2. If term == currentTerm, votedFor is null or candidateId, and candidate's log is at least as complete as local log, grant vote and reset election timeout

## AppendEntries RPC

Invoked by leader to replicate log entries and discover inconsistencies; also used as heartbeat.

### Arguments:

<b>term</b>	leader's term
<b>leaderId</b>	so follower can redirect clients
<b>prevLogIndex</b>	index of log entry immediately preceding new ones
<b>prevLogTerm</b>	term of prevLogIndex entry
<b>entries[]</b>	log entries to store (empty for heartbeat)
<b>commitIndex</b>	last entry known to be committed

### Results:

<b>term</b>	currentTerm, for leader to update itself
<b>success</b>	true if follower contained entry matching prevLogIndex and prevLogTerm

### Implementation:

1. Return if term < currentTerm
2. If term > currentTerm, currentTerm  $\leftarrow$  term
3. If candidate or leader, step down
4. Reset election timeout
5. Return failure if log doesn't contain an entry at prevLogIndex whose term matches prevLogTerm
6. If existing entries conflict with new entries, delete all existing entries starting with first conflicting entry
7. Append any new entries not already in the log
8. Advance state machine with newly committed entries

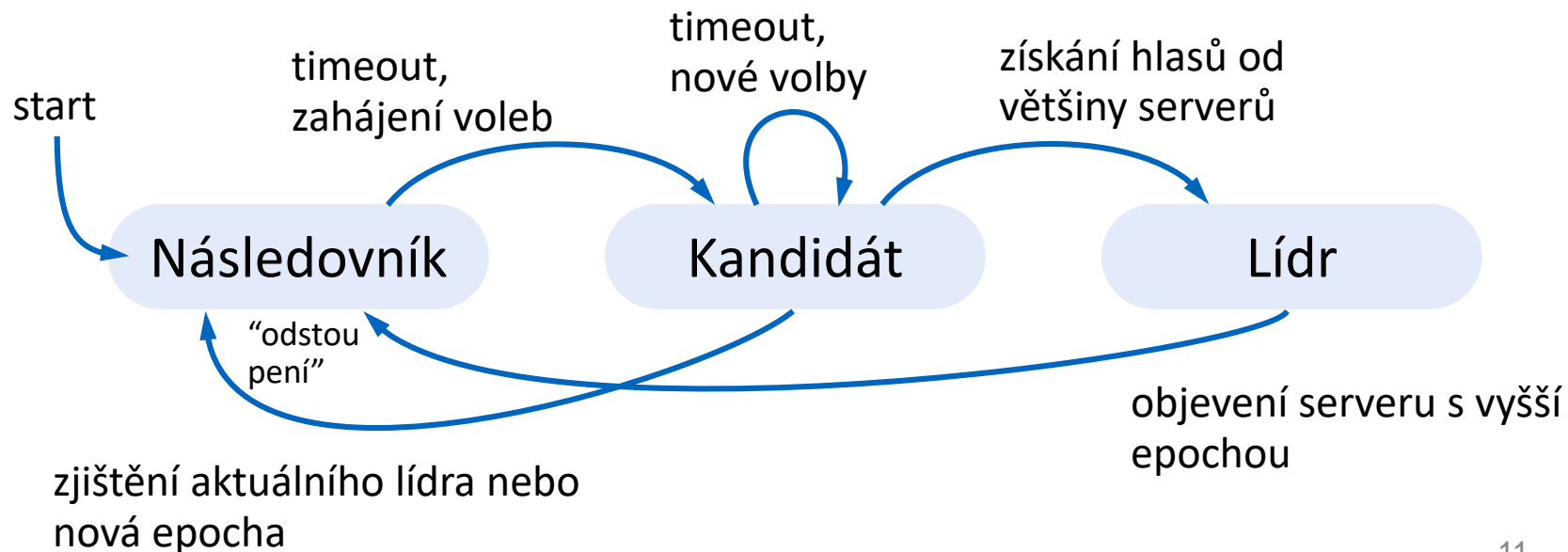
# Stavy serveru

Servery začínají jako **Následovníci**.

- Následovníci očekávají zprávy od Lídra nebo Kandidátů

Lídři posílají **heartbeats** (prázdné zprávy **AppendEntries**), aby si udrželi autoritu.

Jakmile Následovník neobdrží zprávu do **volebního timeoutu** (typicky 100-500ms), předpokládá, že Lídr havaroval a **iniciuje volbu** nového lídra.



# Spuštění voleb

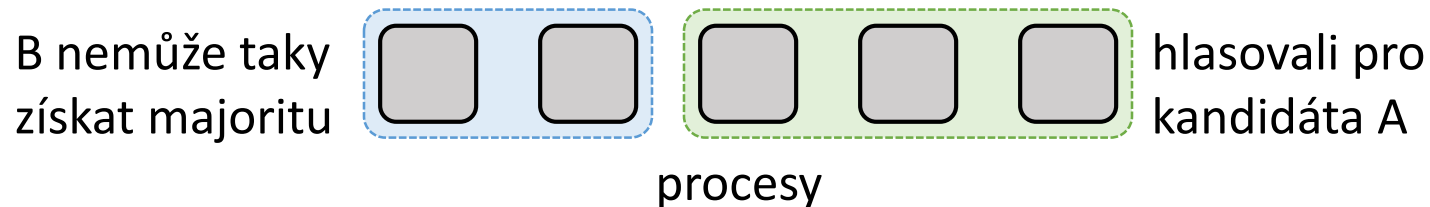
Server, který vyvolá volby, provede následující:

1. Zvýší číslo epochy.
2. Změní svůj stav na KANDIDÁT
3. Zahlasuje pro sebe
4. Pošle **RequestVote** všem ostatním serverům a čeká dokud nenastane jedno z následujících:
  1. Obdrží hlasy od většiny serverů (hladký průběh):
    - Změní stav na LÍDR
    - Pošle **AppendEntries** heartbeats všem ostatním procesům
  2. Přijme zprávu od validního LÍDRA (byl zvolen dříve):
    - Vráť se do stavu NÁSLEDOVNÍK
  3. Nikdo nevyhraje volby (vyprší volební timeout):
    - Zvýší epochu a začne nové volby

# Klíčové vlastnosti voleb

**Bezpečnost:** maximálně jeden vítěz v každé epoše

- každý proces hlasuje pouze jednou v jedné epoše (a hlas persistuje)
- dva kandidáti nemohou získat většinu v jedné epoše



**Živost:** jeden z kandidátů musím časem vyhrát

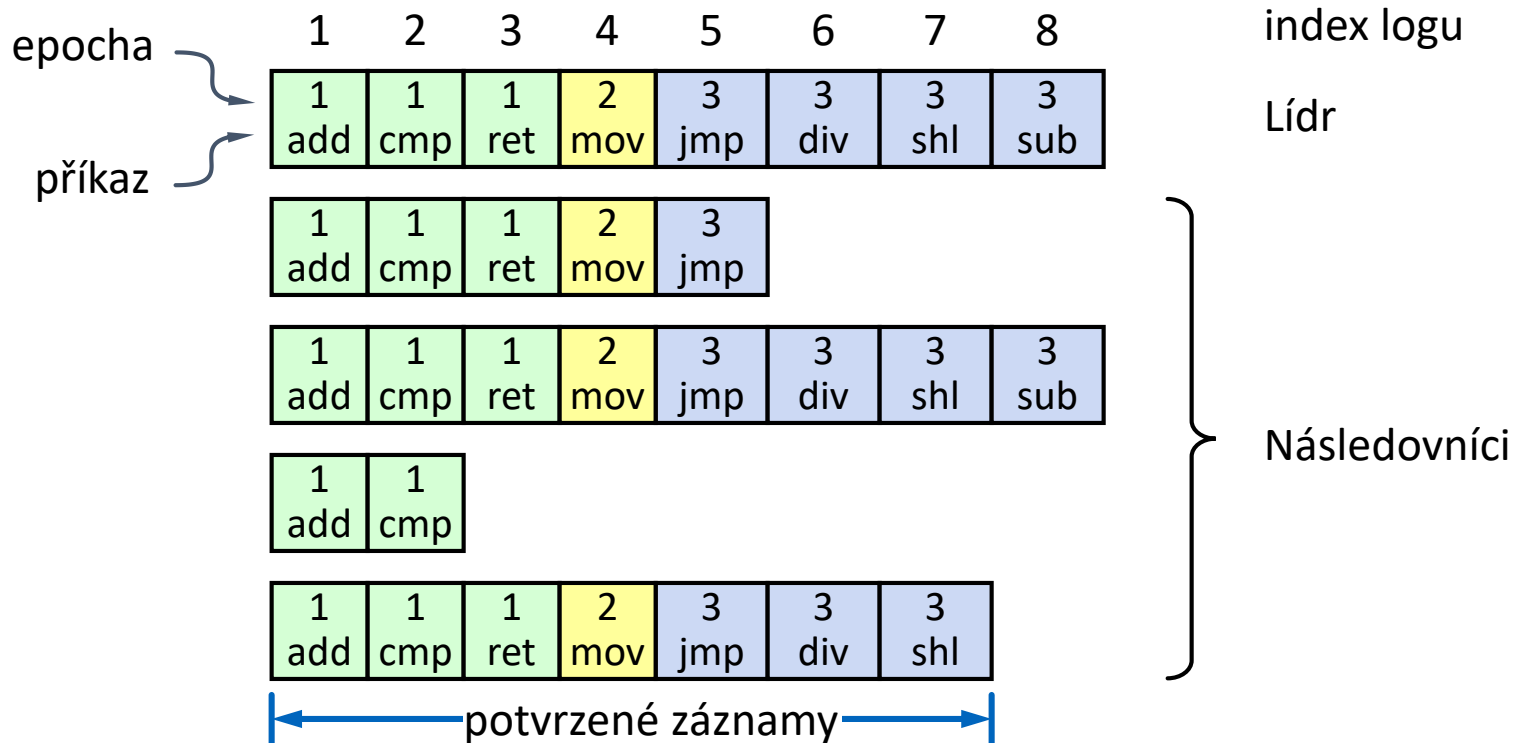
- Každý proces volí volební timeout **náhodně** v intervalu  $[T, 2T]$
- Jeden proces typicky iniciuje volby a zvítězí dříve, než ostatní začnou (volby skončí typicky za ms až desítky ms)
- funguje dobře pokud  $T \gg RTT$  (čas oběhu zpráv) a pokud  $T \ll MTBF$  (střední doba mezi selháními)

<https://raft.github.io/>



# Běžný chod

# Struktura logů



**Záznam** v logu = < index, epocha vzniku, příkaz >

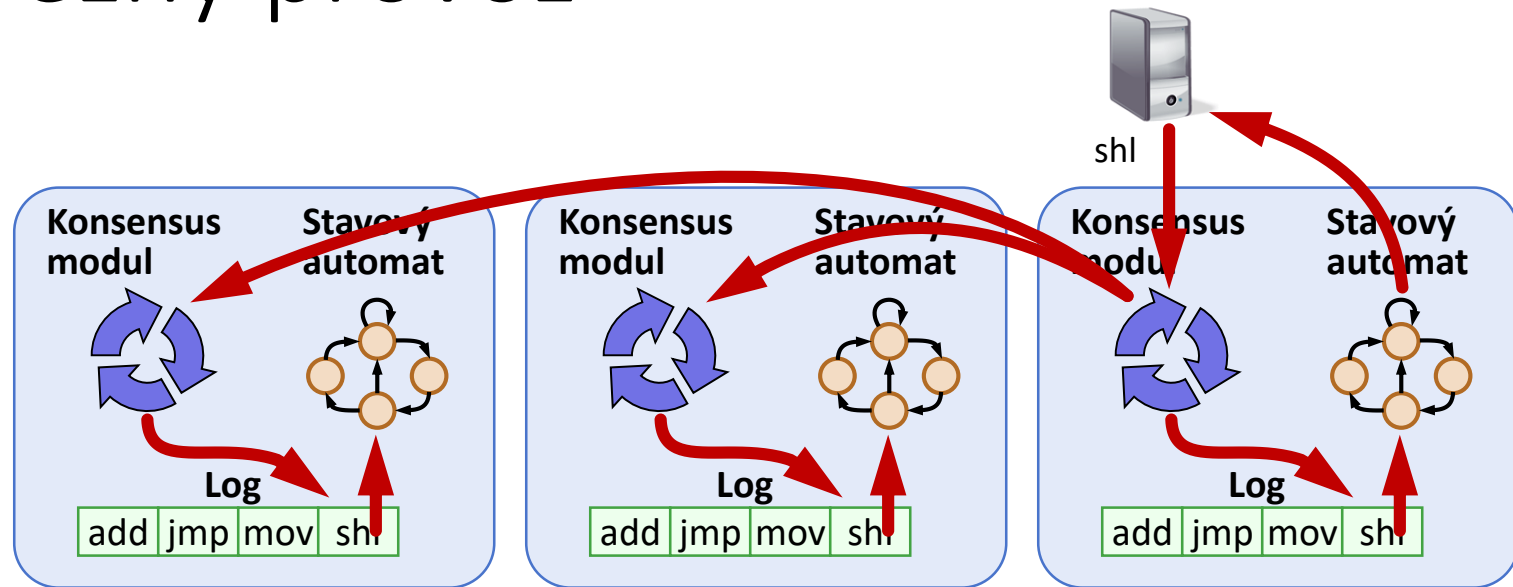
Logy jsou uloženy v **perzistentním uložišti** (disk); tj. přežijí **havárie**

Záznam je **potvrzený (committed)**, je-li známo, že je uložen ve většině procesů

- trvalý: nebude změněn a bude nakonec vykonán stavovým automatem

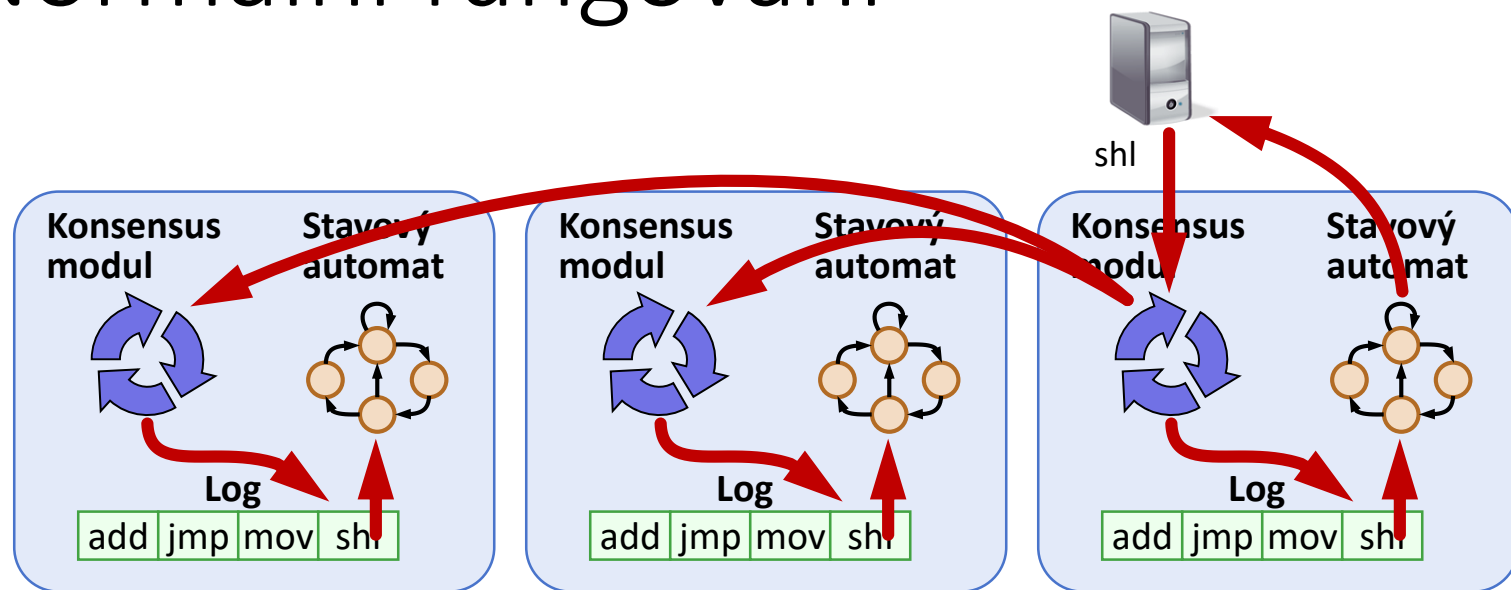


# Běžný provoz



1. Klient pošle příkaz lídrovi.
2. Lídr přidá příkaz na konec svého logu.
3. Lídr pošle zprávu **AppendEntries** následovníkům, typicky paralelně, a čeká na odpovědi.
4. Jakmile je nový záznam potvrzený (committed)
  - Lídr předá příkaz k vykonání svému stavovému automatu a výsledek pošle klientovi.
  - Lídr přidá informaci o potvrzení (commit) do následující zprávy **AppendEntries** pro následovníky
  - Následovníci předají příkaz svým stavovým automatům.

# Normální fungování



Havarování / pomalí následovnicí?

- Lídr opakovaně posílá zprávu **AppendEntries**, dokud doručení neuspěje

V běžném provozu velmi efektivní:

- Stačí úspěšné doručení **AppendEntries** většině procesů a lze odpovědět klientovi

# Konzistence logů

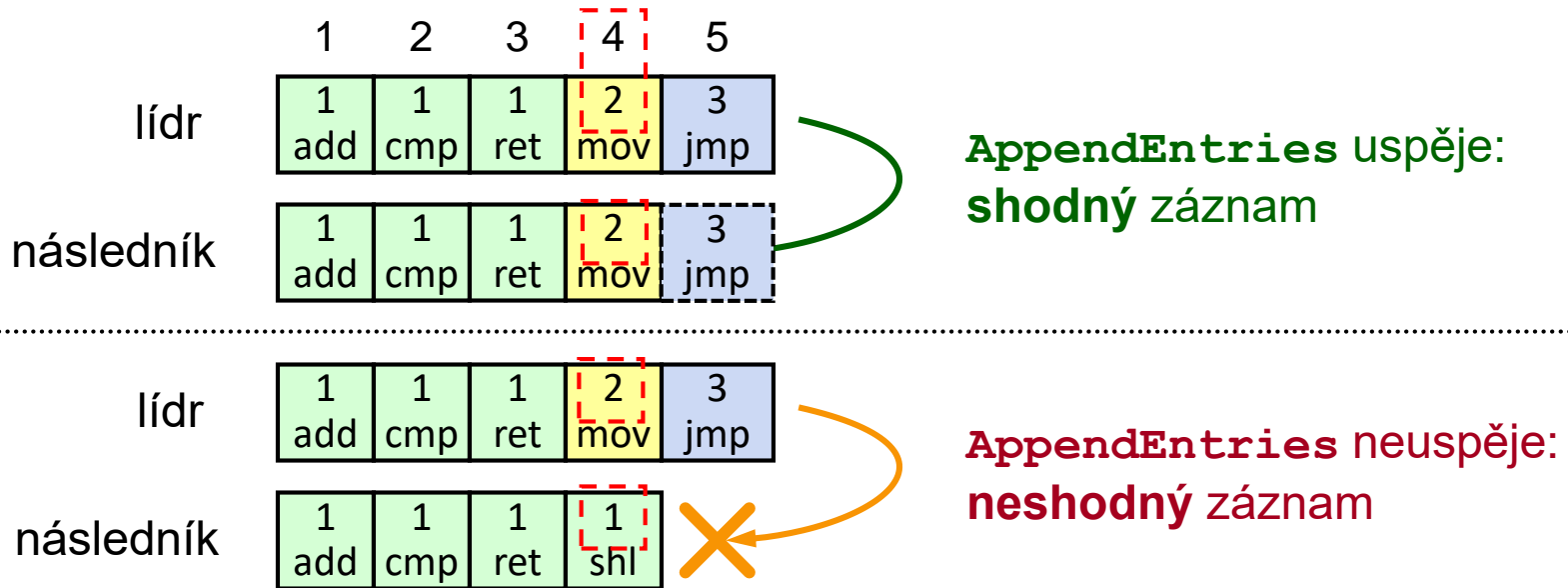
	1	2	3	4	5	6
server 1	1 add	1 cmp	1 ret	2 mov	3 jmp	3 div
server 2	1 add	1 cmp	1 ret	2 mov	3 jmp	4 sub

Raft je navržen tak, aby vynucoval/garantoval následující **invarianty**<sup>1</sup>:

1. Mají-li záznamy logů uložené na různých serverech stejný index a epochu, pak
  - obsahují **stejný příkaz**
  - logy jsou **identické** ve všech **předcházejících** záznamech
2. Je-li daný záznam potvrzený, jsou **potvrzené** i všechny **předcházející** záznamy

<sup>1</sup>Invariant = vlastnost splněna po celou dobu běhu algoritmu

# Kontrola konzistence



**AppendEntries** obsahuje  $\langle \text{index}, \text{term} \rangle$  záznamu předcházejícího nově přidávané záznamy.

Následovník musí obsahovat **shodný záznam**; jinak je zápis odmítnut.

Kontrola shodnosti předcházejícího záznamu implementuje **indukční krok** a zajišťuje konzistenci logu.



# Změna lídra

# Konsistence logů

Během **normální** fungování zůstávají logy lídra a následovníků **konzistentní**.

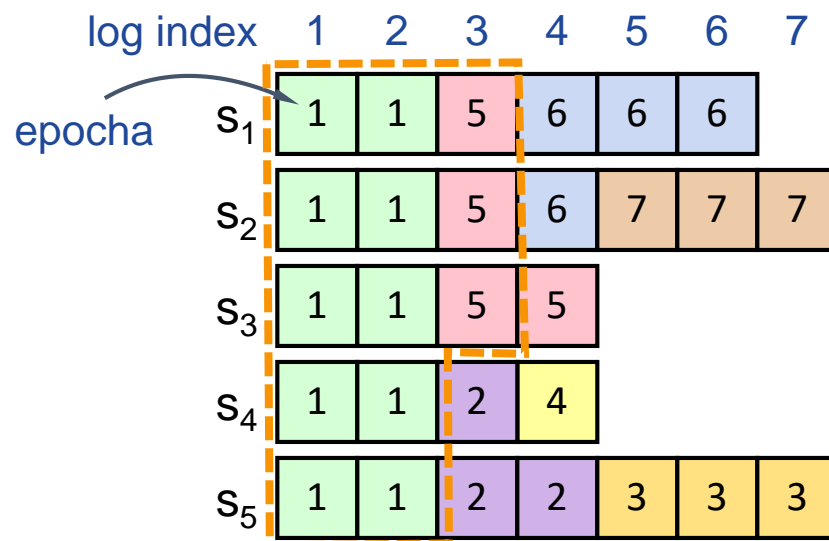
Pokud ovšem lídr **havaruje** a je zvolen nový lídr, mohou být logy lídra a následovníků **nekonzistentní**.

# Změny lídra

Log nového lídra vždy reprezentuje „pravdu“ (**správné záznamy**) – po jeho zvolení není potřeba žádné speciální kroky, vykoná logiku běžného chodu.

- logika běžného chodu *časem* udělá logy následovníků identické s logem nového lídra
- záznamy logu předchozího lídra mohou být částečně replikovány (ale nepotvrzeny) – budou postupně eliminovány

Dojde-li k několika haváriím lídrů po sobě, může být v ložích jednotlivých procesů řada přebytečných záznamů, které budou postupně eliminovány.



# Bezpečnost

Obecně nutná bezpečnostní garance pro replikaci

Jakmile je příkaz ze záznamu logu vykonán některým stavovým automatem, nesmí žádný jiný stavový automat vykonat *jiný* příkaz pro stejný záznam.

**Bezpečnostní invariant Raftu:** Jakmile lídr prohlásí záznam v logu za potvrzený, jakýkoliv budoucí lídr bude mít tento záznam ve svém logu.

Invariant Raftu implikuje bezpečnostní garanci:

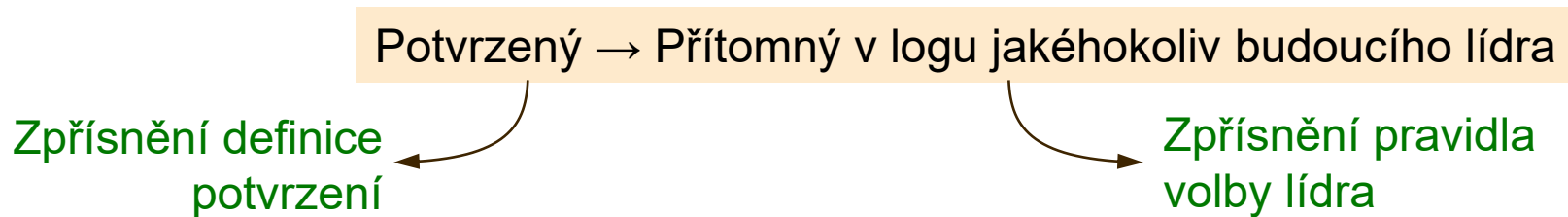
- lídři **nikdy nepřepisují** záznamy ve svých lozích (pouze přidává)
- pouze záznamy **v logu lídra** mohou být **potvrzeny**
- záznamy (příkazy) musí být v logu **potvrzeny předtím**, než jsou **vykonány** stavovým automatem



# Zpřísnění Raftu

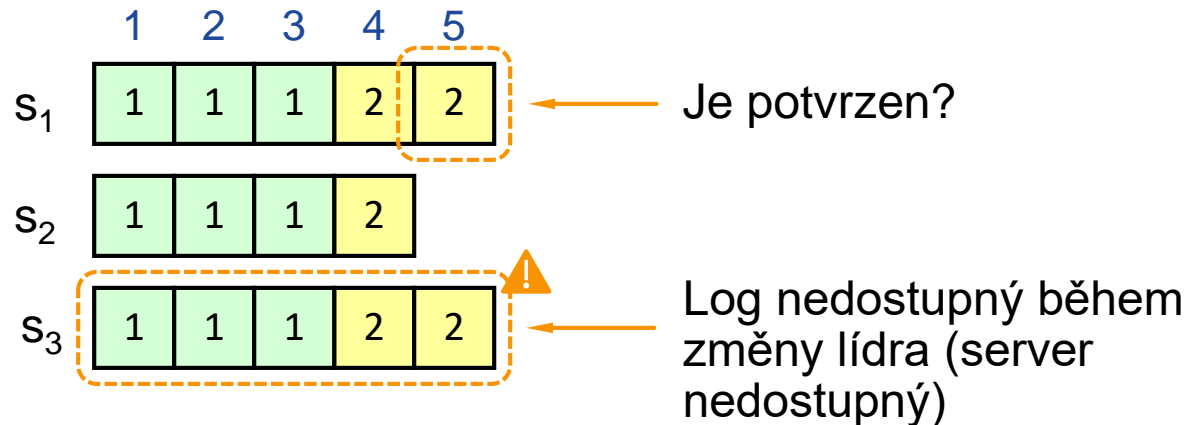
Dosavadní logika fungování Raftu bezpečnostní invariant **negarantuje!**

→ nutno zpřísnit pravidla:



# Zpřísnění výběru lídra

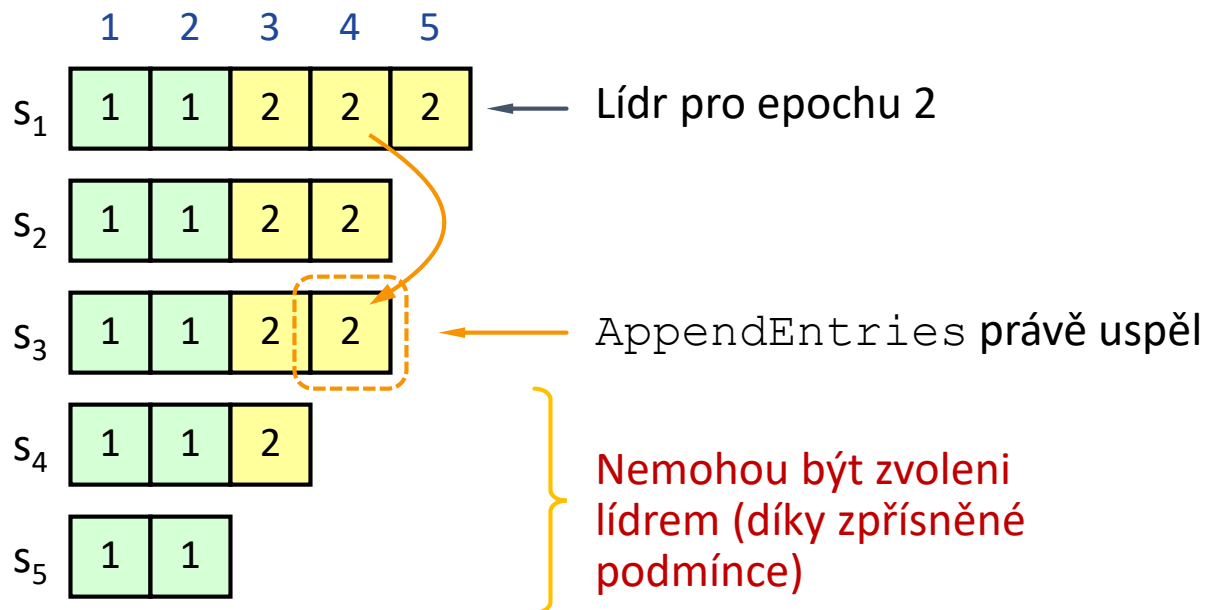
Nelze říct, které záznamy byly potvrzeny



Raft volí kandidáta, který má nejúplnější log. Kandidáti do zprávy **RequestVote** vloží index a epochu posledního záznamu svého logu

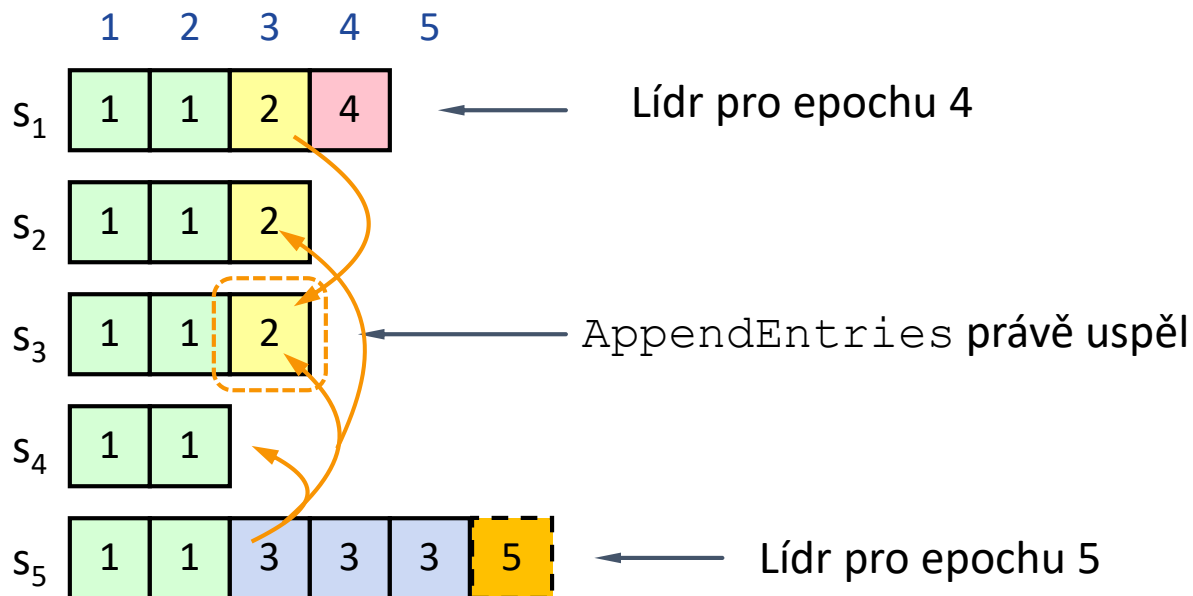
- Volící server hlas pro kandidáta odmítne, pokud jeho vlastní log je **úplnější**, tj. pokud má na konci záznam s vyšší epochou nebo stejnou epochou, ale vyšším indexem.
- Lídr tedy bude mít **nejúplnější log** mezi většinou procesů, kterou byl zvolen.

# Potvrzování záznamu z aktuální epochy



Záznam 4 je bezpečně potvrzen: jakýkoliv lídr pro epochu tři musí obsahovat v logu záznam 4.

# Potvrzování záznamu z dřívější epochy



## Záznam 3 **není bezpečně** potvrzený

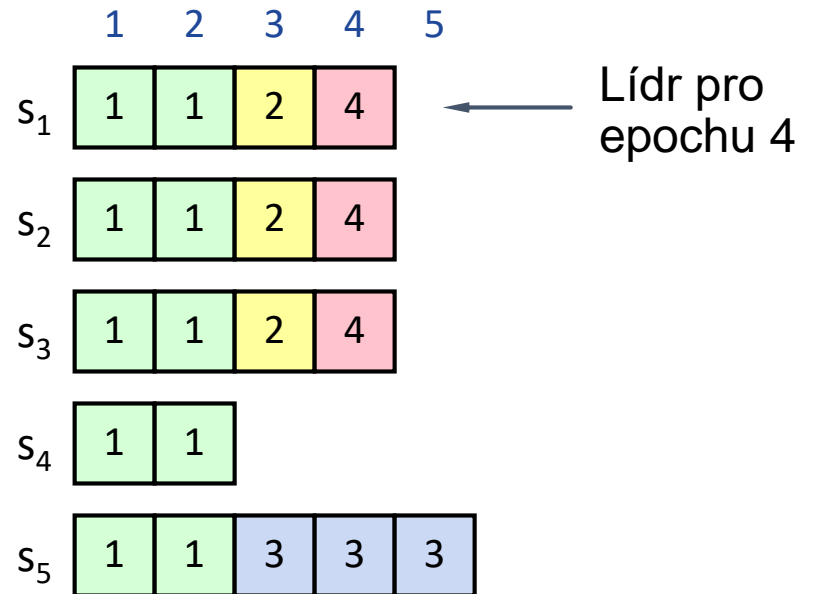
- $S_5$  může být zvolen jako lídr pro epochu 5
- Byl by-li zvolen, přepíše záznam 3 v  $S_1, S_2, S_3$

# Nová pravidla pro potvrzování

Aby lídr považoval záznam za potvrzený:

1. záznam musí být uložený na většině serverů
2. **aspoň jeden nový záznam** z lídrovy aktuální epochy musí být taky na většině serverů

Příklad: Jakmile je záznam 4 potvrzen,  $S_5$  nemůže být zvolen lídrem pro epochu 5 a záznamy 3 a 4 jsou bezpečně potvrzeny.

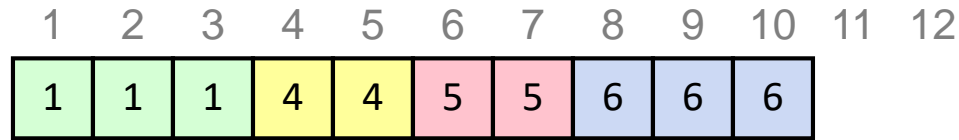


Kombinace nové pravidla pro výběr lídra a zpřísněné definice potvrzování **garantuje bezpečnostní invariant** Raftu

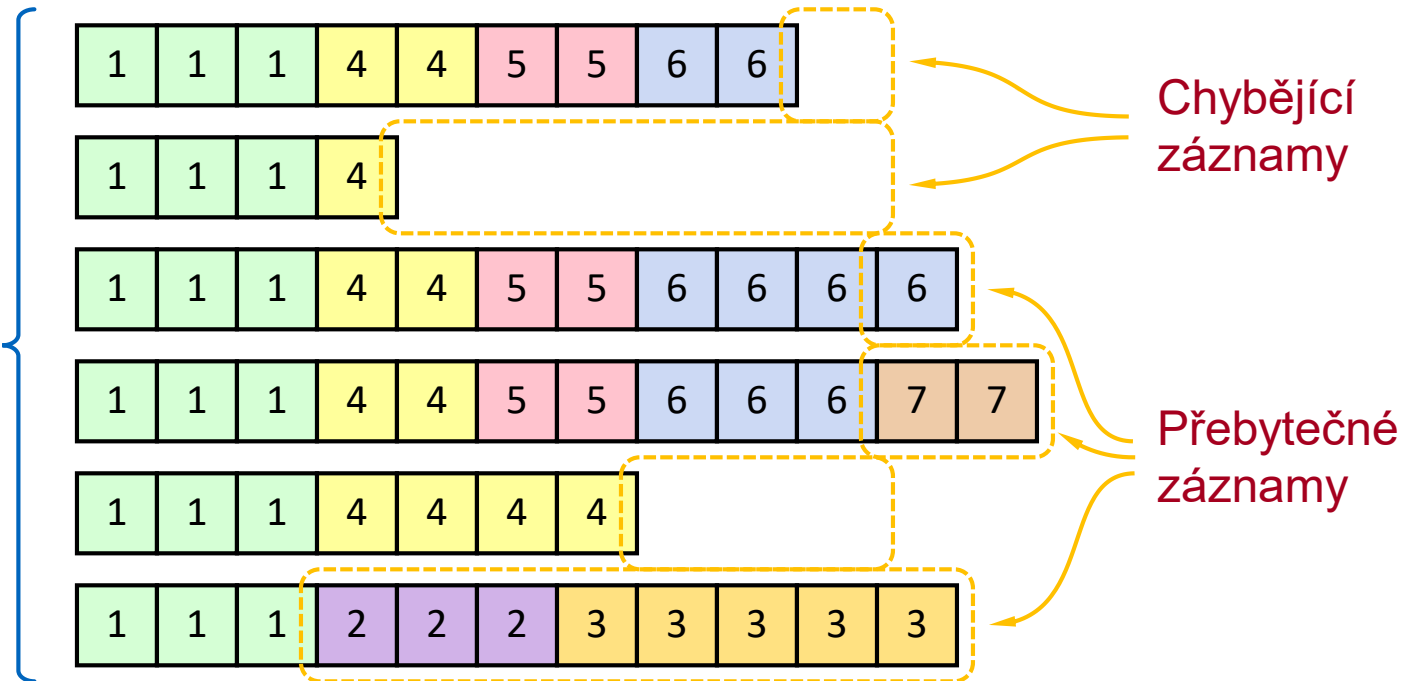
# Komplikace: nekonzistence logu

Změna lídra mohou vést k nekonzistencím logu

Lídr pro epochu 8



Možní  
následovníci



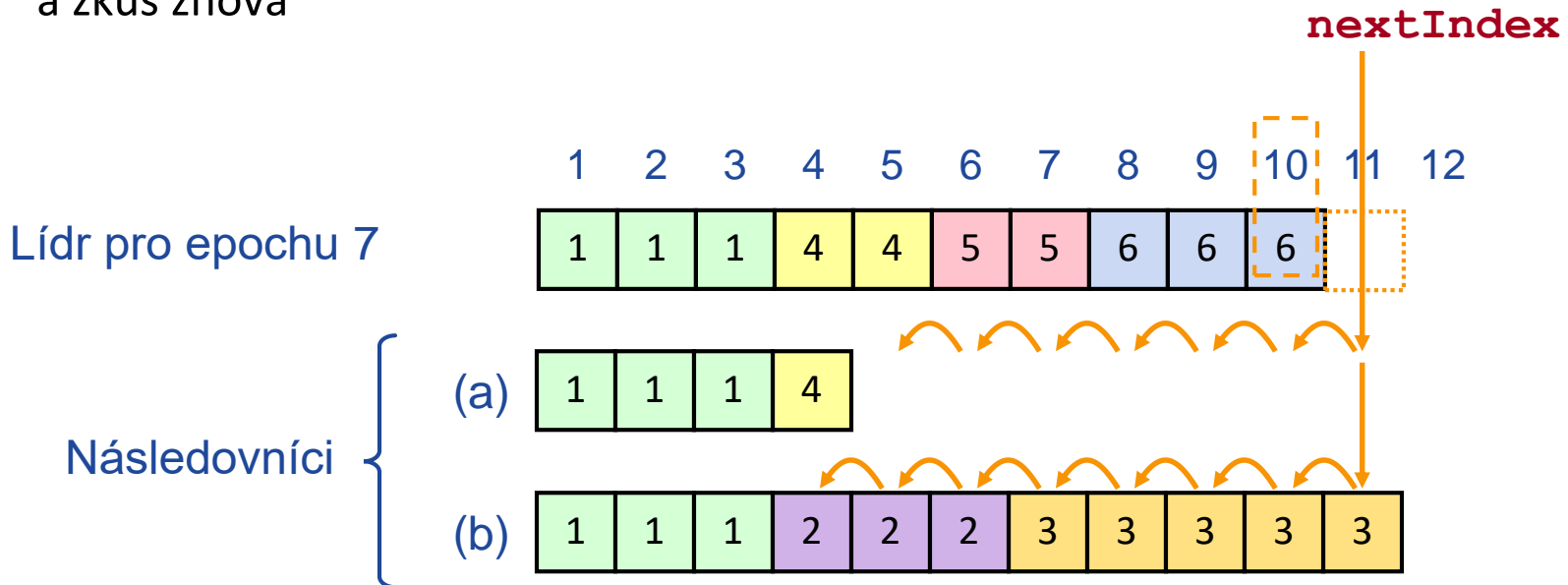
# Oprava logů následovníků

Nový lídr musí udělat logy následovníků konzistentní se svým logem, tj. smazat přebytečné záznamy a doplnit chybějící záznam.

Lídr udržuje proměnou **nextIndex** pro *každého* následovníka:

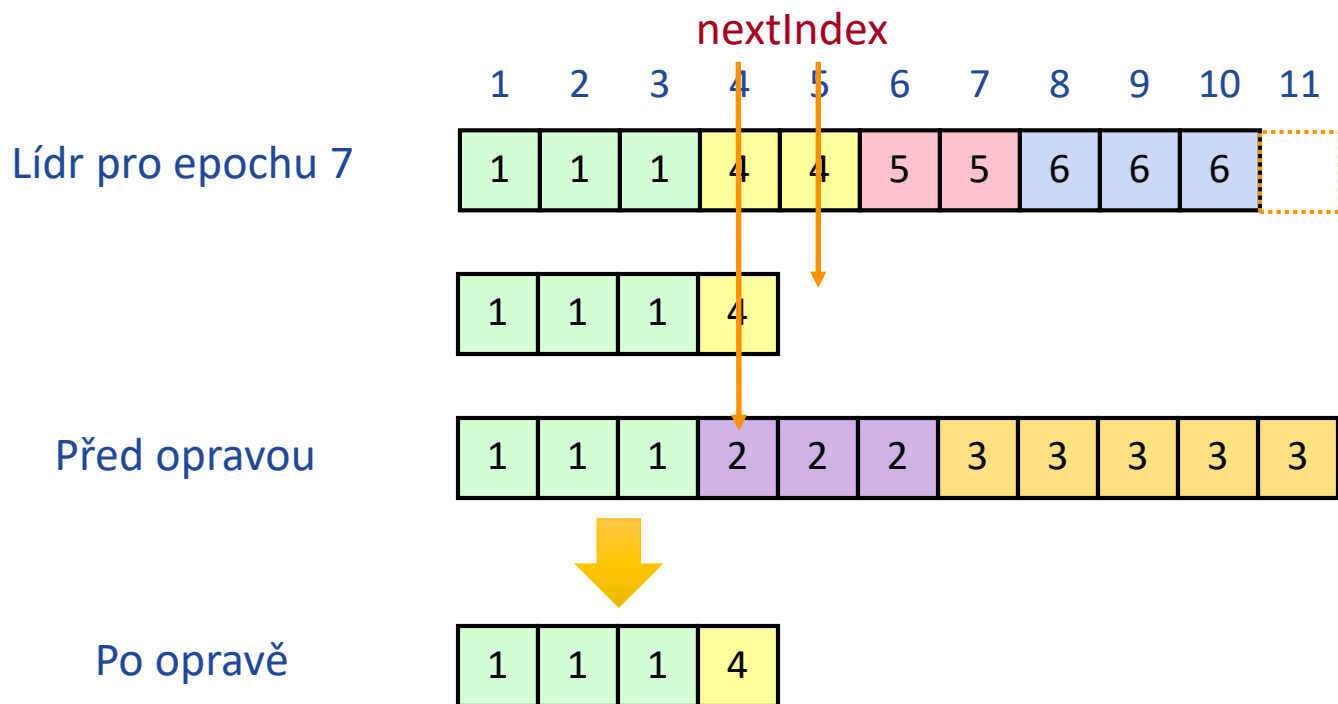
- index další záznamu logu, který by měl být odeslán následníkovi
- inicializován na 1 + index poslední záznamu lídra

Pokud kontrola konzistence **AppendEntries** selže, s niž nextIndex o jednu a zkus znova



# Oprava logů následovníků

Pokud následovník přepíše nekonzistentní záznam, odstraní i **všechny následující** záznamy.







# Neutralizace starého lídra

# Neutralizace starých lídrů

Sesazený lídr **nemusí** být **trvale** havarovaný

- přechodné odpojení od sítě
- jiné procesy zvolí nového lídra
- starý lídr se znovu připojí a pokusí se potvrdit svoje záznamy

Epochy slouží k **detekci neaktuálního** lídra

- každá zpráva obsahuje epochu odesílatele
- je-li epocha odesílatele starší, zpráva je odmítnuta, odesílatel se změní na Následovníka a aktualizuje si epochu
- je-li epocha příjemce starší, tak příjemce se změní na Následovníka, aktualizuje si epochu a následně zprávu normálně zpracuje

**Volby aktualizují epochy** většiny serverů → sesazení lídři nemohou potvrdit nové záznamy



# Klientský protokol

# Protokol klienta

Klienti posílají příkazy lídrovi

- Není-li lídr známý, kontaktují libovolný server a ten je případně přesměruje na lídra

Lídr pouze vrací odezvu na příkaz poté, co je příkaz **zalogován**, **potvrzen** a následně vykonán **lídrem**.

Pokud **nepřijde** v časovém limitu **odezva** na požadavek (např. lídr havaroval):

- klient vybere (náhodně) jiný server
- a po případném přesměrování nakonec odešle příkaz novému lídrovi

# Protokol klienta: jediné vykonání

Lídr může havarovat poté, co vykonal příkaz, ale před odesláním odpovědi

→ Riziko **opakovaného** vykonání příkazu.

Řešení: Pro zajištění právě **jednoho vykonání** příkazu klient vloží unikátní ID příkazu do každého požadavku

- Toto ID je uložené v záznamech v logu

Před přijetím požadavku lídr zkontroluje, zda-li už nemá záznam se stejným ID ve svém logu

- Pokud ne → příkaz vykoná;
- Pokud ano → příkaz odmítne.

# Souhrn

Problém konsensu je v jádru mnoha problémů v DS.

V asynchronním DS **nelze** při přítomnosti selhání **konsensus vyřešit** ve smyslu bezpečnosti a živosti.

Praktická řešení garantují **bezpečnost**.

Raft je moderní algoritmus pro replikaci logů / výpočtů.

Je využíván v řadě reálných DS.

Literatura:

- [Raft] Ongaro, D. and Ousterhout, J.K., 2014, June. In search of an understandable consensus algorithm. In *USENIX Annual Technical Conference*. [\[link\]](#)