

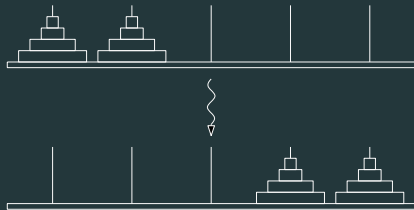
Čas a detekce selhání v distribuovaných systémech

B4B36PDV – Paralelní a distribuované výpočty

- Opakování z minulého cvičení
- Detekce selhání v distribuovaných systémech
- Čas a uspořádání událostí v distribuovaných systémech
- Zadání 6. domácí úlohy

Odevzdání semestrální
práce se blíží!

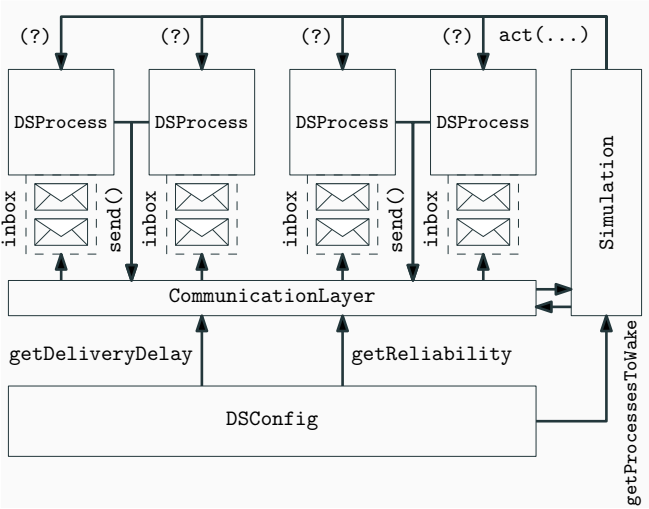
Neděle 7.5. 23:59 CET



Opakování z minulého cvičení

<http://goo.gl/a6BEMb>

DSand framework



Jakými všemi způsoby může být následující kód proveden?

```
// CONFIG:
    getProcessesToWake() { return {"1", "2", "3"}; }
// PROCESS:
    int time = 0;
    int nid = Integer.parseInt(id);
    public void act() {
        time++;
        if(nid == time && nid != 3)
            send("3", new DummyMessage());
        while(!inbox.isEmpty()){
            Message m = inbox.poll();
            System.out.println(m.sender);
        }
    }
```

Detekce selhání v DS

I have a feeling
that something
went wrong...

ZZZ...

by cicakkia '07

I have a feeling
that something
went wrong...

ZZZ...

by cicakkia '07

Vzpomeňte si na BFS

Co když spící/mrtvá hlava
leží na nejkratší cestě?

Když nám „umře“ důležitý proces, musíme být schopní se se situací vypořádat. (Jinak nám celý DS zhavaruje)

Když nám „umře“ důležitý proces, musíme být schopní se se situací vypořádat. (Jinak nám celý DS zhavaruje)

⚠ Ale to musíme nejdřív zjistit, že „umřel“!

Když nám „umře“ důležitý proces, musíme být schopní se se situací vypořádat. (Jinak nám celý DS zhavaruje)

⚠ Ale to musíme nejdřív zjistit, že „umřel“!

Jak na to?

- Heartbeats jsou odesílány periodicky (každých T „vteřin“)
- Nedostane-li proces heartbeat od procesu p_j po dobu $T + \tau$ „vteřin“, považuje p_j za mrtvý

- Heartbeats jsou odesílány periodicky (každých T „vteřin“)
- Nedostane-li proces heartbeat od procesu p_j po dobu $T + \tau$ „vteřin“, považuje p_j za mrtvý
- Centralizovaný heartbeat
- Kruhový heartbeat
- Všichni-všem (all-to-all) heartbeating

- Heartbeats jsou odesílány periodicky (každých T „vteřin“)
 - Nedostane-li proces heartbeat od procesu p_j po dobu $T + \tau$ „vteřin“, považuje p_j za mrtvý
 - Centralizovaný heartbeat
 - Kruhový heartbeat
 - Všichni-všem (all-to-all) heartbeating
-

Doprogramujte detekování selhání procesu na základě (all-to-all) heartbeating

Doimplementujte logiku detekování selhání procesu na základě (all-to-all) heartbeating v `DetectorProcess.java`. Následně spusťte scénář `MainFD.java`, ve kterém máte zajištěno, že selže právě jeden proces.

- **Úplnost:** každé selhání je časem detekováno aspoň jedním funkčním procesem

- **Úplnost:** každé selhání je časem detekováno aspoň jedním funkčním procesem
- **Přesnost:** nedochází k označení funkčního procesu za havarovaný

- **Úplnost:** každé selhání je časem detekováno aspoň jedním funkčním procesem
- **Přesnost:** nedochází k označení funkčního procesu za havarovaný
- **Rychlost:** čas do okamžiku, kdy první proces detekuje selhání

- **Úplnost:** každé selhání je časem detekováno aspoň jedním funkčním procesem
- **Přesnost:** nedochází k označení funkčního procesu za havarovaný
- **Rychlost:** čas do okamžiku, kdy první proces detekuje selhání
- **Škálovatelnost:** ani při velkém počtu agentů nedojde k zahlcení systému

- **Úplnost:** každé selhání je časem detekováno aspoň jedním funkčním procesem
- **Přesnost:** nedochází k označení funkčního procesu za havarovaný
- **Rychlost:** čas do okamžiku, kdy první proces detekuje selhání
- **Škálovatelnost:** ani při velkém počtu agentů nedojde k zahlcení systému

All-to-all heartbeating:

:-(Přesnost

:-(Škálovatelnost

Čas a uspořádání událostí v DS (1. část)

V centralizovaném systému je čas konzistentní...
(procesy typicky sdílí jediné hodiny)

Sdílené hodiny můžeme snadno využít pro:

- Koordinaci
(„výpočet zahájíme v 11:47:23“)
- Uspořádání kroků výpočtu
(logování, uspořádání procesů při přístupu do kritické sekce atd.)
- ... a jiné

V distribuovaném systému máme problém...
(každý proces má vlastní hodiny)

Možné problémy které hodiny mohou mít:

V distribuovaném systému máme problém...
(každý proces má vlastní hodiny)

Možné problémy které hodiny mohou mít:

- "Clock skew"
(hodiny A: 16:24, hodiny B: 16:32)

V distribuovaném systému máme problém...
(každý proces má vlastní hodiny)

Možné problémy které hodiny mohou mít:

- "Clock skew"
(hodiny A: 16:24, hodiny B: 16:32)
- "Clock drift"
(Hodiny začnou ve stejný čas, např. 12:00, ale po jedné hodině jedny ukazují 12:59 a druhé 13:01)

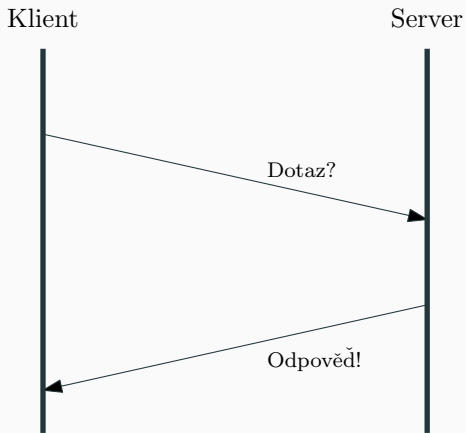
V distribuovaném systému máme problém...
(každý proces má vlastní hodiny)

Možné problémy které hodiny mohou mít:

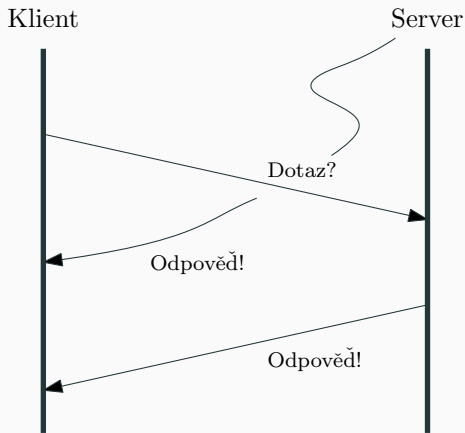
- "Clock skew"
(hodiny A: 16:24, hodiny B: 16:32)
- "Clock drift"
(Hodiny začnou ve stejný čas, např. 12:00, ale po jedné hodině jedny ukazují 12:59 a druhé 13:01)

⚠ Často ale nepotřebujeme znát přesný čas, stačí znát skutečnou souslednost událostí.

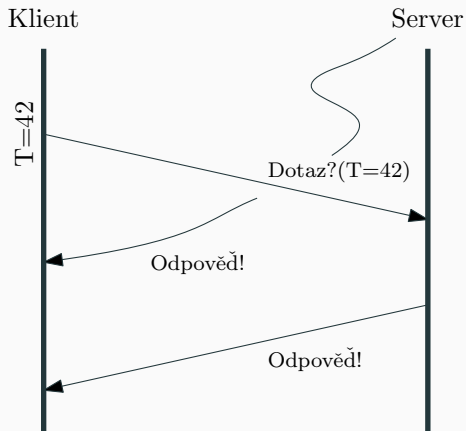
Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



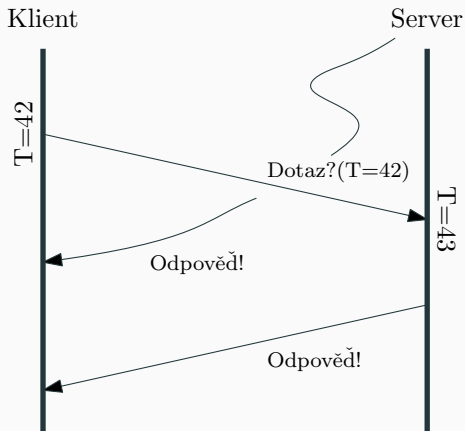
Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



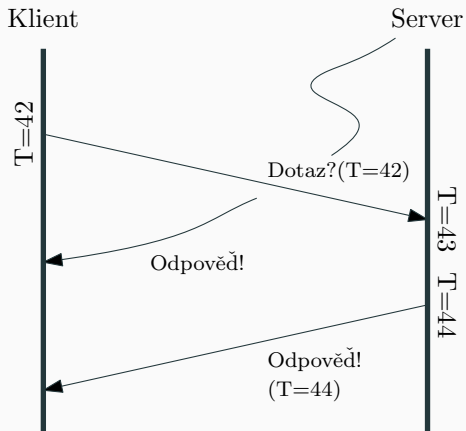
Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



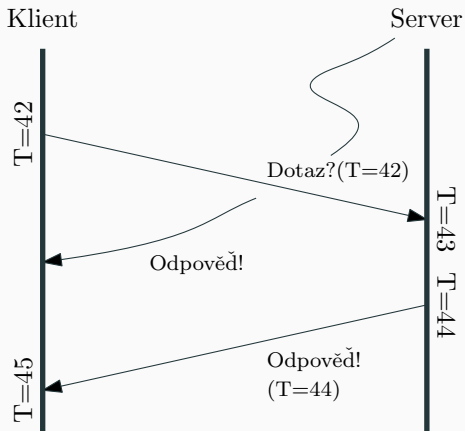
Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



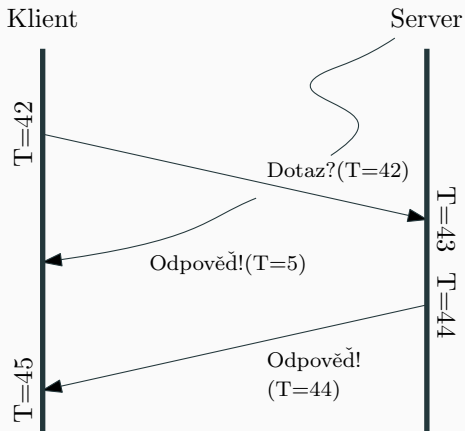
Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



Vztah příčiny a následku je v DS klíčový!
(například, odpověď na dotaz následuje až po položení dotazu)



Právě jsme v našem DS zavedli logický čas! :-)
(konkrétně Lamportovy skalární hodiny)

Právě jsme v našem DS zavedli logický čas! :-)
(konkrétně Lamportovy skalární hodiny)

Logický čas splňuje pouze kauzalitu!

- Každé události e přiřadíme časovou značku $T(e)$
- Pokud je událost e příčinou události e' , pak platí $T(e) < T(e')$
(Ne nutně ale naopak!)

Lamportův algoritmus

Lamportův algoritmus

1. Každý proces má svoje lokální logické hodiny

```
int logicalTime = 0
```

Lamportův algoritmus

1. Každý proces má svoje lokální logické hodiny

```
int logicalTime = 0
```

2. Před každou významnou událostí (obzvlášť posláním zprávy!) si proces lokální čas posune

```
++logicalTime
```


Lamportův algoritmus

1. Každý proces má svoje lokální logické hodiny

```
int logicalTime = 0
```

2. Před každou významnou událostí (obzvlášť posláním zprávy!) si proces lokální čas posune

```
++logicalTime
```

3. Každé zprávě přiřadíme časovou značku `msg.T = logicalTime`
(Tím říkáme přijímajícímu procesu, ať si upraví svůj čas!)

Lamportův algoritmus

1. Každý proces má svoje lokální logické hodiny
`int logicalTime = 0`
2. Před každou významnou událostí (obzvlášť posláním zprávy!) si proces lokální čas posune
`++logicalTime`
3. Každé zprávě přiřadíme časovou značku `msg.T = logicalTime`
(Tím říkáme přijímajícímu procesu, ať si upraví svůj čas!)
4. Přijetí zprávy je následkem jejího odeslání – pak musí platit $T(e) < T(e')$
Po přijetí zprávy `msg` si proto musíme zaktualizovat svůj `logicalTime`:

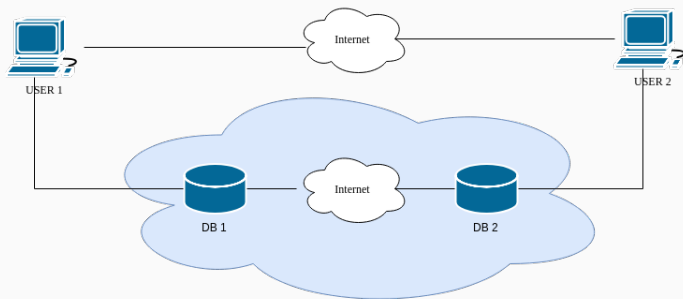
$$\text{logicalTime} = 1 + \max\{\text{logicalTime}, \text{msg.T}\}$$

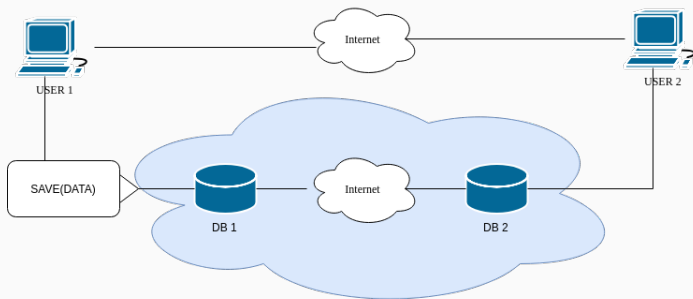
Lamportův algoritmus

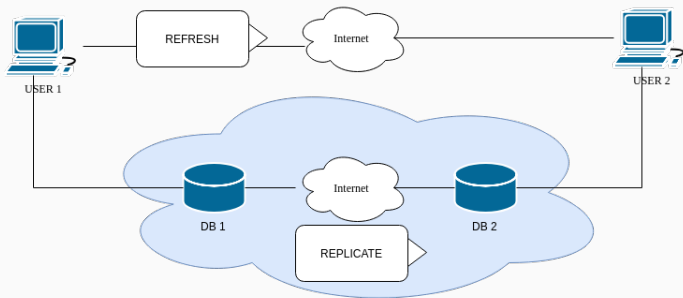
1. Každý proces má svoje lokální logické hodiny
`int logicalTime = 0`
2. Před každou významnou událostí (obzvláště posláním zprávy!) si proces lokální čas posune
`++logicalTime`
3. Každé zprávě přiřadíme časovou značku `msg.T = logicalTime`
(Tím říkáme přijímajícímu procesu, ať si upraví svůj čas!)
4. Přijetí zprávy je následkem jejího odeslání – pak musí platit $T(e) < T(e')$
Po přijetí zprávy `msg` si proto musíme zaktualizovat svůj `logicalTime`:

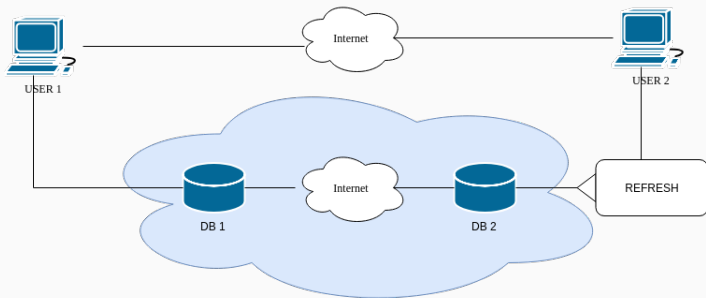
$$\text{logicalTime} = 1 + \max\{\text{logicalTime}, \text{msg.T}\}$$

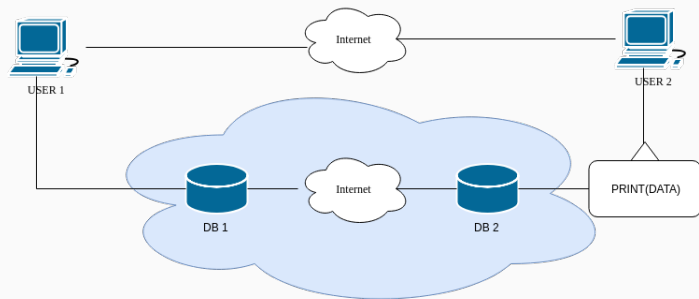
⚠ Skalární hodiny jsou stavebním kamenem mnoha algoritmů v DS!

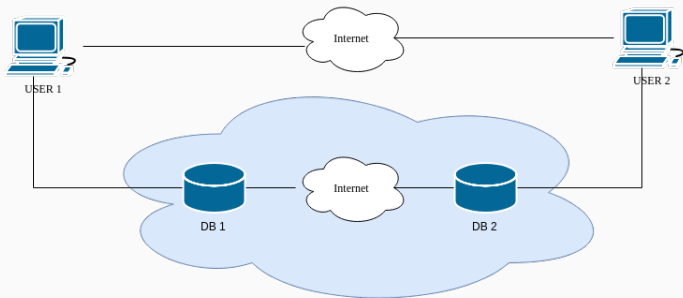












K čemu bychom zde mohli chtít používat logické hodiny?

Doprogramujte Lamportovy logické hodiny

Doimplementujte logiku Lamportových logických hodin ve třídě `ScalarClock.java`. Následně spusťte scénář `ScalarClockRun.java`.

Doprogramujte Lamportovy logické hodiny

Doimplementujte logiku Lamportových logických hodin ve třídě `ScalarClock.java`. Následně spusťte scénář `ScalarClockRun.java`.

Co je v našem systému špatně?

Doprogramujte Lamportovy logické hodiny

Doimplementujte logiku Lamportových logických hodin ve třídě `ScalarClock.java`. Následně spusťte scénář `ScalarClockRun.java`.

Co je v našem systému špatně?

⚠ Replikace může být pomalá. Druhý klient tak může číst stará data!

Doprogramujte Lamportovy logické hodiny

Doimplementujte logiku Lamportových logických hodin ve třídě `ScalarClock.java`. Následně spusťte scénář `ScalarClockRun.java`.

Co je v našem systému špatně?

⚠ Replikace může být pomalá. Druhý klient tak může číst stará data!

Jsme to schopní detekovat skalárními hodinami?

Doimplementujte metodu `isCausalityForProcessViolated`

Pak zkuste spustit scénář `ScalarDSConfigBombarding`

Jak protokol upravit, aby nedocházelo k porušení
kauzality?

Jak protokol upravit, aby nedocházelo k porušení
kauzality?

Možností je mnoho, například:

Jak protokol upravit, aby nedocházelo k porušení kauzality?

Možností je mnoho, například:

- Před odesláním **REFRESH** zprávy si počkat na potvrzení od databáze (Odeslání **REFRESH** zprávy je kauzálním následkem úspěšné replikace)

Jak protokol upravit, aby nedocházelo k porušení kauzality?

Možností je mnoho, například:

- Před odesláním **REFRESH** zprávy si počkat na potvrzení od databáze (Odeslání **REFRESH** zprávy je kauzálním následkem úspěšné replikace)
- Pozdržet vyhodnocení dotazu do doby, než replikace proběhne (Druhému uživateli můžeme poslat, že má požadovat data zapsaná nejdříve v daném logickém čase)

Jak protokol upravit, aby nedocházelo k porušení kauzality?

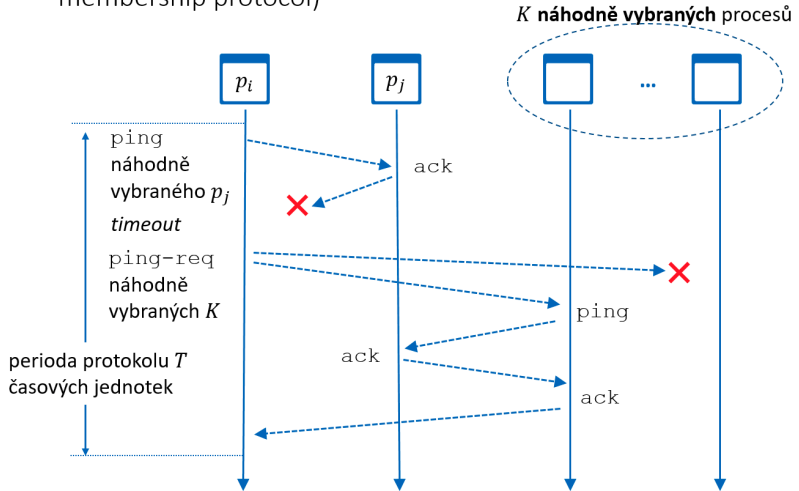
Možností je mnoho, například:

- Před odesláním **REFRESH** zprávy si počkat na potvrzení od databáze (Odeslání **REFRESH** zprávy je kauzálním následkem úspěšné replikace)
 - Pozdržet vyhodnocení dotazu do doby, než replikace proběhne (Druhému uživateli můžeme poslat, že má požadovat data zapsaná nejdříve v daném logickém čase)
- ⚠** Obecně chceme, aby události e_1 , e_2 , které mají proběhnout po sobě (tj. například čtení až po replikaci) byly ve vztahu kauzální závislosti.

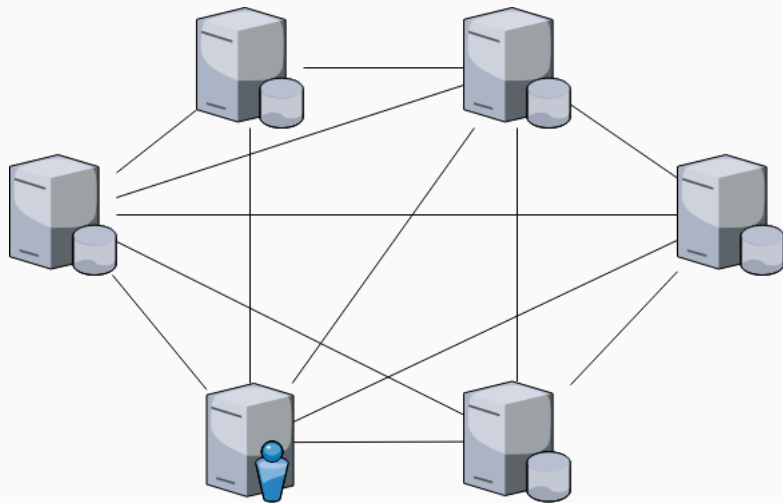
Zadání domácí úlohy

SWIM Failure Detector

(Scalable weakly consistent infection-style process group membership protocol)



Distribuovaná databáze na komoditních serverech



Naimplementujte SWIM detekci selhání a zajistěte, že

1. zbytečně nevytěžuje síť;
2. detekuje všechny “mrtvé” procesy s rozumnou rychlostí; a
3. je dostatečně přesné.

Zpracování musí být **distribuované**, procesy si nesahají vzájemně do paměti!

Díky za pozornost!

Budeme rádi za Vaši
zpětnou vazbu! →



[https://forms.gle/
yi7FWBEw3mxgnJ9P7](https://forms.gle/yi7FWBEw3mxgnJ9P7)