

Как при помощи бумаги, карандаша и алгоритма Raft достичь консенсуса

Ярослав Дынников

Picodata



HighLoad ++

Слайды: <https://rosik.github.io/2023-highload>

О чем речь

Кластер — это группа процессов, работающих совместно и представляющихся пользователю единым компьютерным ресурсом.

О чем речь

Кластер — это группа процессов, работающих совместно и представляющихся пользователю единым компьютерным ресурсом.

Задача

- Есть несколько серверов.
- Надо достичь консенсуса.

О чем речь

Кластер — это группа процессов, работающих совместно и представляющихся пользователю единым компьютерным ресурсом.

Задача

- Есть несколько серверов.
- Надо достичь консенсуса.
- В ненадежной сети.

О чем речь

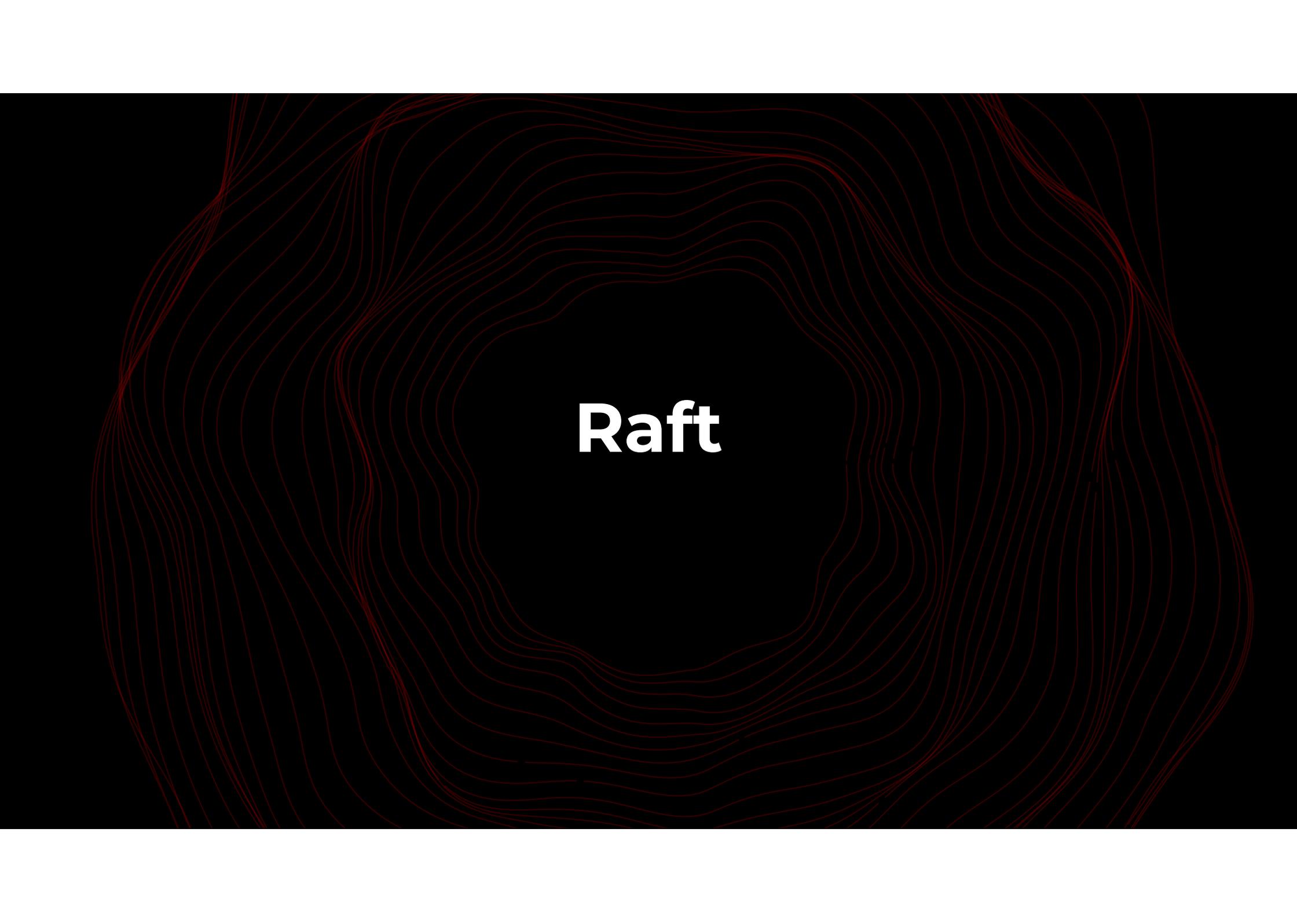
Кластер — это группа процессов, работающих совместно и представляющихся пользователю единым компьютерным ресурсом.

Задача

- Есть несколько серверов.
- Надо достичь консенсуса.
- В ненадежной сети.

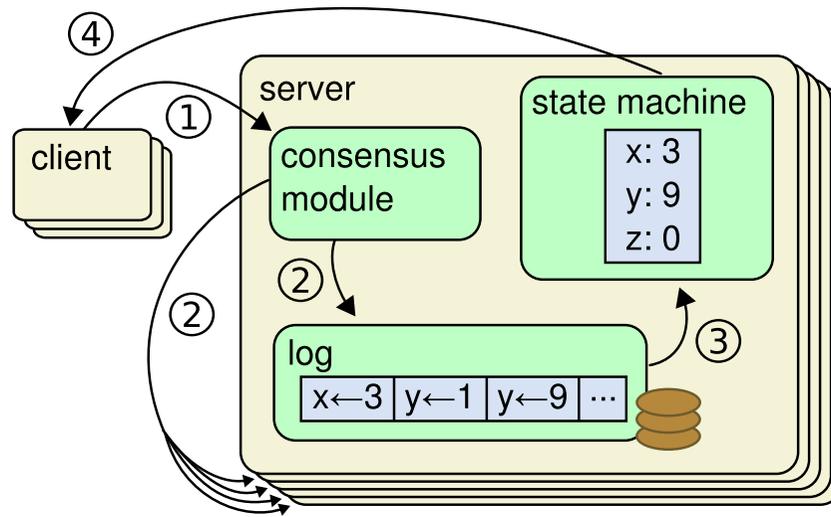
Решение — Raft

- In search of Understandable Consensus Algorithm.
- Diego Ongaro and John Ousterhout. Stanford University.
- <https://raft.github.io>

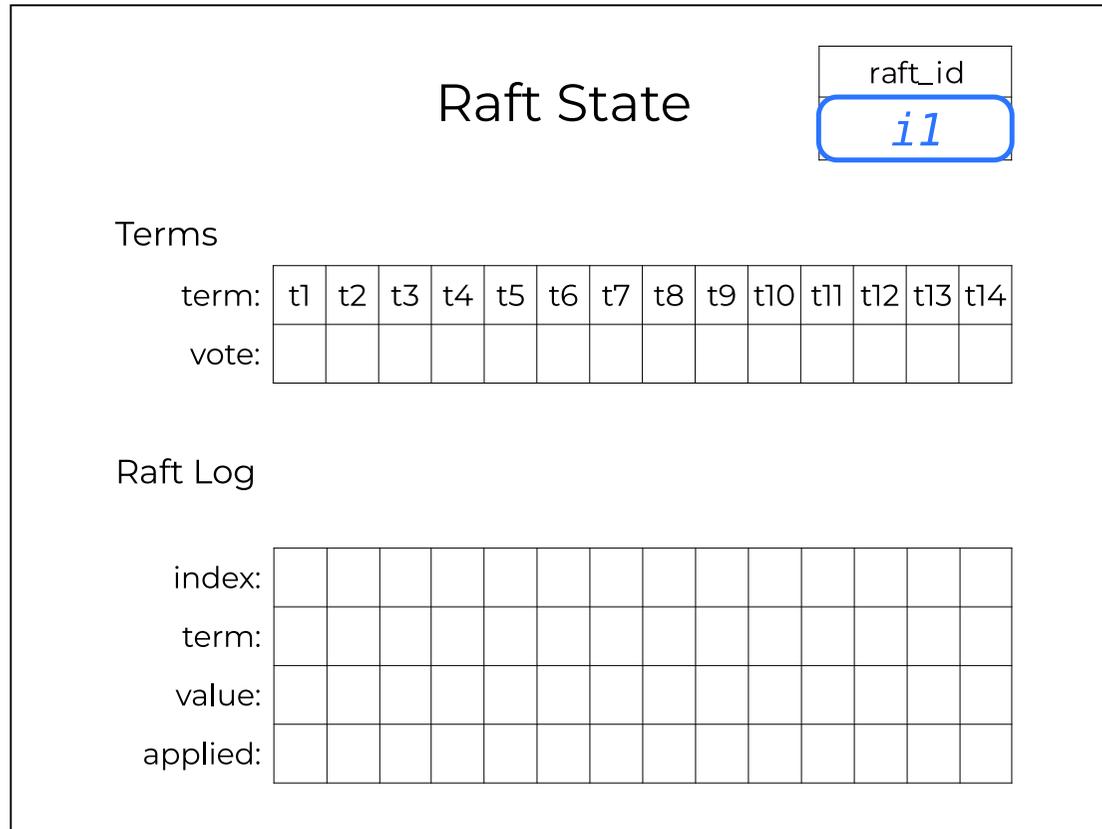
The background of the image is a solid black field. Overlaid on this field are numerous thin, wavy lines in a dark red or maroon color. These lines are arranged in a pattern that resembles concentric, irregular waves or ripples, creating a sense of movement and depth. The lines are most densely packed in the center and become more sparse towards the edges.

Raft

Реплицируемый конечный автомат

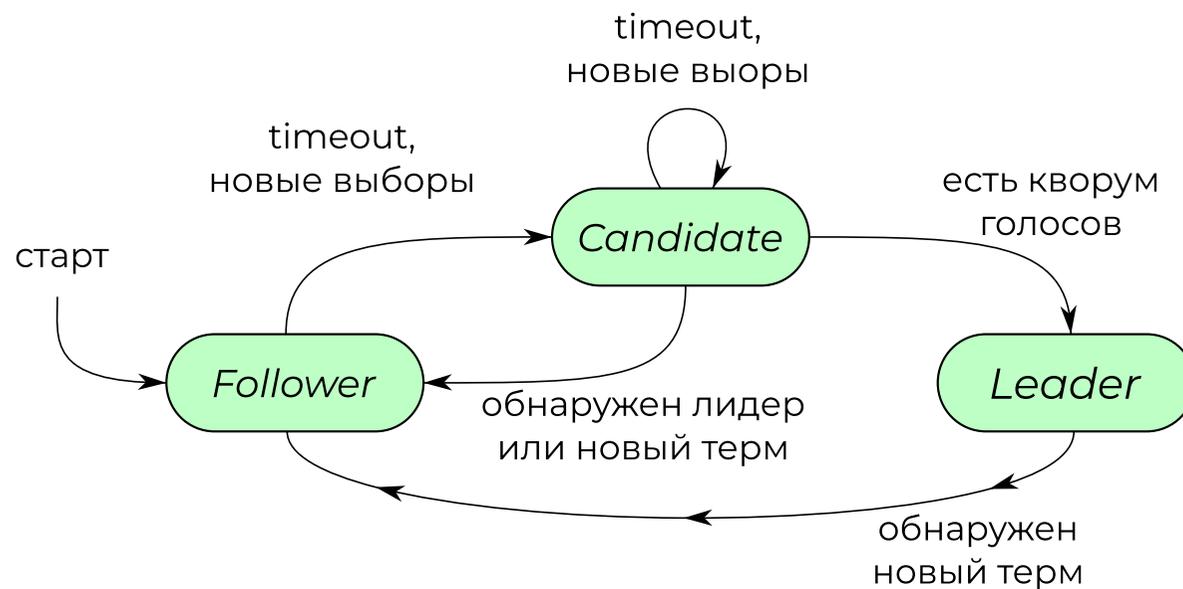


Персистентное хранилище



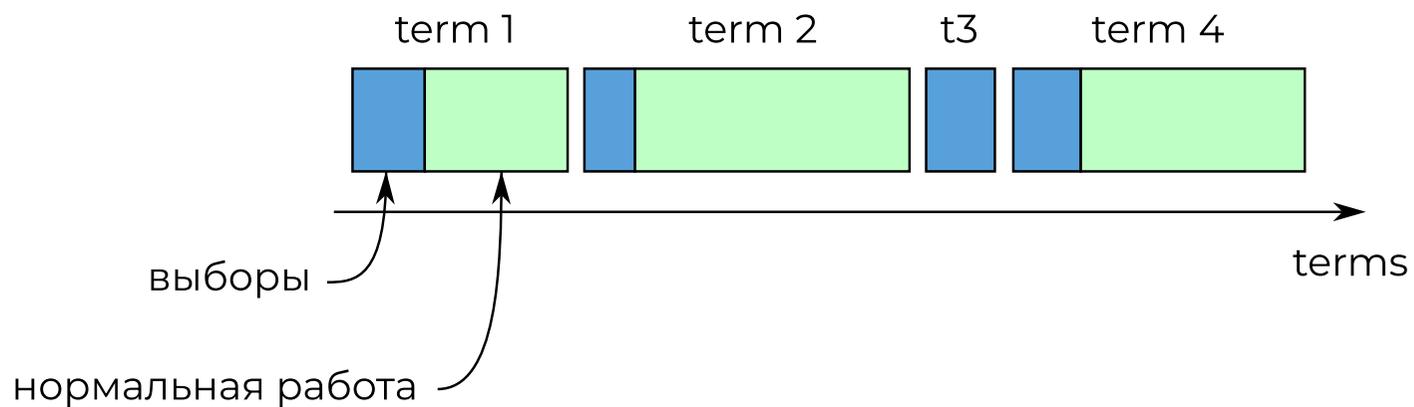
Лидер, фолловер, кандидат

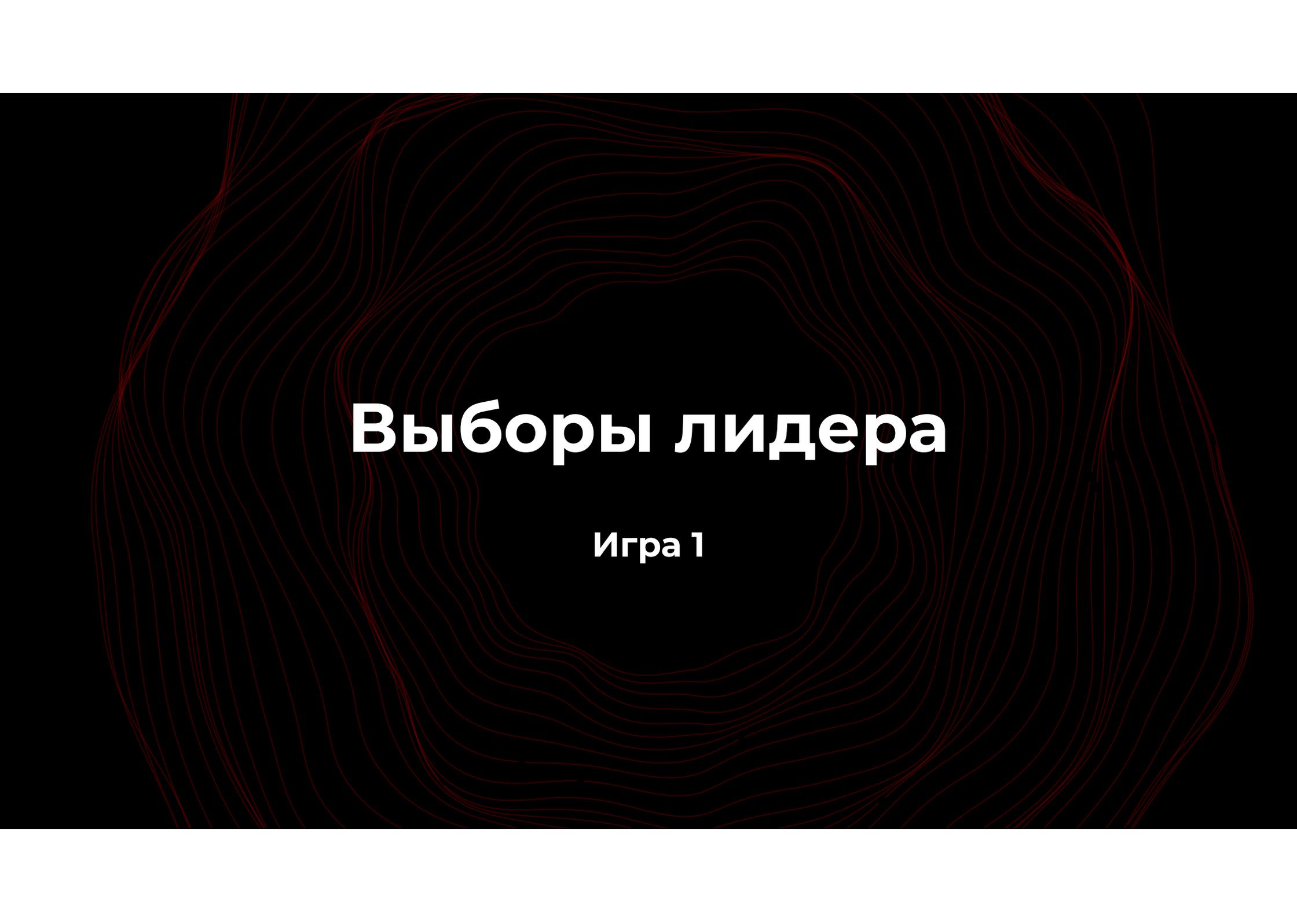
- *Leader* единственный пишет в журнал + пингует окружающих
- *Follower* пассивен, не отправляет никаких запросов
- *Candidate* проводит голосование



Термы

Терм — это отрезок времени неопределенной длины. Он начинается с выборов, после которых единственный лидер управляет кластером.





Выборы лидера

Игра 1

i1, начинайте выборы



Raft State

raft_id
<i>i1</i>

Terms

term:	<i>t1</i>	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:														
term:														
value:														
applied:														

RequestVoteRPC

term:	<i>t1</i>
candidateId:	<i>i1</i>
lastLogIndex:	<i>0</i>
lastLogTerm:	<i>t0</i>

Response

	<i>i1</i>	i2	i3	i4	i5	i6	i7	i8
ok	<input checked="" type="checkbox"/>							

i2-i8: «OK»



RequestVoteRPC

term:	<i>t1</i>
candidateId:	<i>i1</i>
lastLogIndex:	<i>0</i>
lastLogTerm:	<i>t0</i>

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓	✓						

Raft State

raft_id
<i>i2</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:														
term:														
value:														
applied:														

i1: «Y-xyy!»

Raft State

raft_id
<i>i1</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:														
term:														
value:														
applied:														

RequestVoteRPC

term:	<i>t1</i>
candidateId:	<i>i1</i>
lastLogIndex:	<i>0</i>
lastLogTerm:	<i>t0</i>

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓	✓	✓	✓	✓	✓	✓	✓



Репликация журнала

Игра 2



іІ, заповняйте raft-журнал

Raft State

raft_id
<i>1</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>1</i>													

Raft Log

index:	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>						
term:	<i>t1</i>	<i>t1</i>	<i>t1</i>						
value:	<i>H</i>	<i>i</i>	<i>g</i>	<i>h</i>	<i>L</i>	<i>o</i>	<i>a</i>	<i>d</i>						
applied:														

i1: «Заперсистьте!»



Raft State

raft_id
<i>i1</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:	1	2	3	4	5	6	7	8						
term:	t1	t1	t1						
value:	H	i	g	h	L	o	a	d						
applied:														

AppendEntriesRPC

term:	<i>t1</i>
leaderId:	<i>i1</i>
index:	<i>1</i> <i>2</i> <i>3</i> <i>4</i>
value:	<i>H</i> <i>i</i> <i>g</i> <i>h</i>
leaderCommit:	<i>0</i>

Response

	<i>i1</i>	i2	i3	i4	i5	i6	i7	i8
ok	<input checked="" type="checkbox"/>							

i2-i8: «OK»



AppendEntriesRPC

term:	<i>t1</i>
leaderId:	<i>i1</i>
index:	<i>1</i> <i>2</i> <i>3</i> <i>4</i>
value:	<i>H</i> <i>i</i> <i>g</i> <i>h</i>
leaderCommit:	<i>0</i>

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓	✓						

Raft State

raft_id
<i>i2</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>										
term:	<i>t1</i>	<i>t1</i>	<i>t1</i>	...										
value:	<i>H</i>	<i>i</i>	<i>g</i>	<i>h</i>										
applied:														

i1: «Отлично!»

Raft State

raft_id
<i>i1</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:	1	2	3	4	5	6	7	8						
term:	t1	t1	t1						
value:	H	i	g	h	L	o	a	d						
applied:	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>										

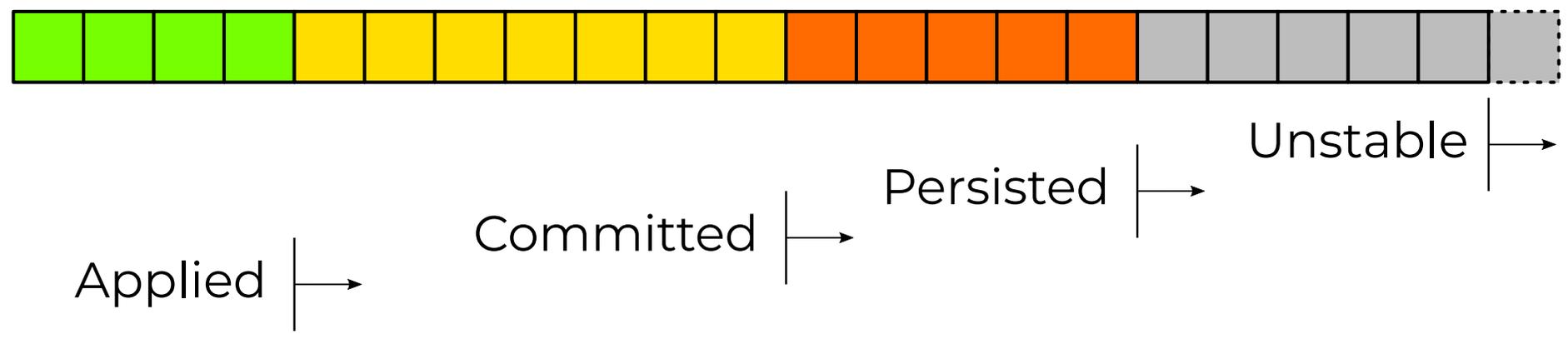
AppendEntriesRPC

term:	<i>t1</i>
leaderId:	<i>i1</i>
index:	1 2 3 4
value:	H i g h
leaderCommit:	<i>0</i>

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	<input checked="" type="checkbox"/>							

Состояние записей



i1: «Реплицируйтесь!»



Raft State

raft_id
<i>i1</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:	1	2	3	4	5	6	7	8						
term:	t1	t1	t1						
value:	H	i	g	h	L	o	a	d						
applied:	✓	✓	✓	✓										

AppendEntriesRPC

term:	<i>t1</i>
leaderId:	<i>i1</i>
index:	5 6 7 8
value:	L o a d
leaderCommit:	4

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓							



i2, i3, сохраняйте записи

AppendEntriesRPC

term:	<i>t1</i>
leaderId:	<i>i1</i>
index:	<i>5</i> <i>6</i> <i>7</i> <i>8</i>
value:	<i>L</i> <i>o</i> <i>a</i> <i>d</i>
leaderCommit:	4

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓	✓						

Raft State

raft_id
<i>i2</i>

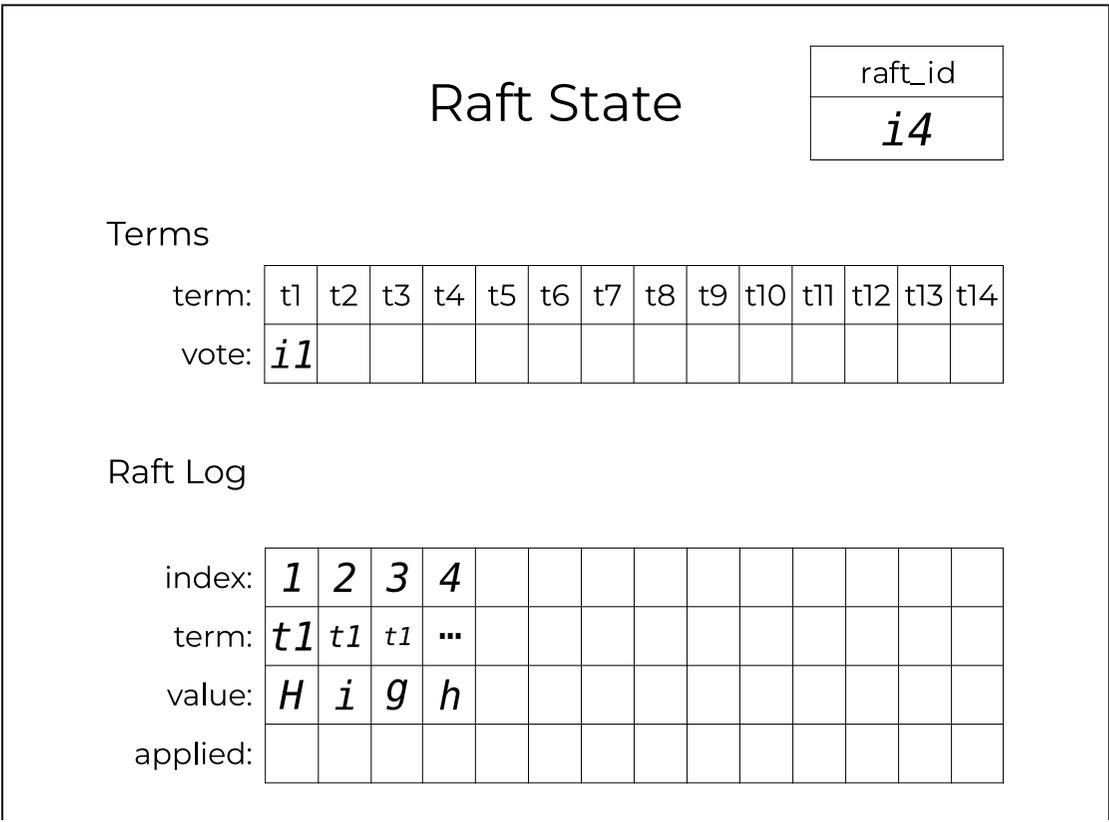
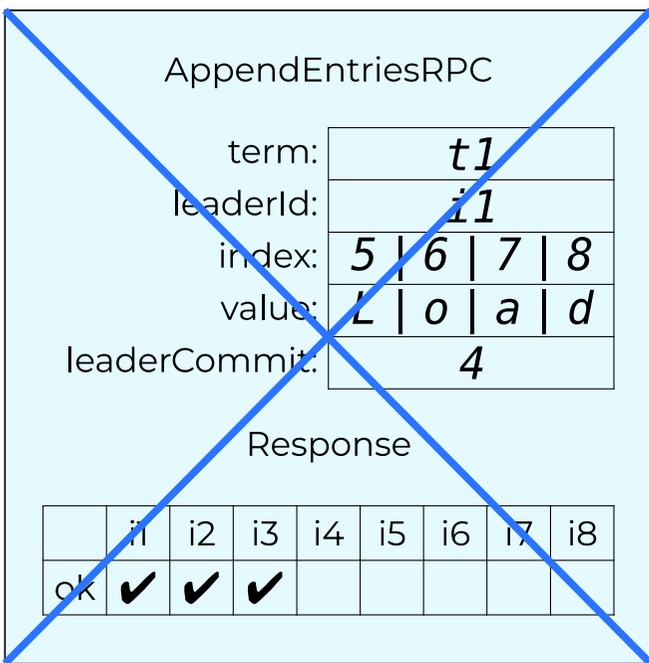
Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>													

Raft Log

index:	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>						
term:	<i>t1</i>	<i>t1</i>	<i>t1</i>							
value:	<i>H</i>	<i>i</i>	<i>g</i>	<i>h</i>	<i>L</i>	<i>o</i>	<i>a</i>	<i>d</i>						
applied:	✓	✓	✓	✓										

i4, "потеряйте" сообщение



Терм без лидера

Игра 3

i4 и i1, вы offline

Переверните ваши листки

i8, начинайте выборы



Raft State

raft_id
<i>i8</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>	<i>i8</i>												

Raft Log

index:	1	2	3	4										
term:	t1	t1	t1	...										
value:	H	i	g	h										
applied:														

RequestVoteRPC

term:	t2
candidateId:	i8
lastLogIndex:	4
lastLogTerm:	t1

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok								✓



i2, i3: «He-a»

RequestVoteRPC

term:	t2
candidateId:	i8
lastLogIndex:	4
lastLogTerm:	t1

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok		X						✓

S may only vote for L if:
 L.lastLogTerm > S.lastLogTerm or
 (L.lastLogTerm == S.lastLogTerm and
 L.lastLogIndex ≥ S.lastLogIndex)

Raft State

raft_id
i2

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	i1	-												

Raft Log

index:	1	2	3	4	5	6	7	8						
term:	t1	t1	t1						
value:	H	i	g	h	L	o	a	d						
applied:	✓	✓	✓	✓										

i8: «😞»



raft_id
<i>i8</i>

Raft State

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>	<i>i8</i>												

Raft Log

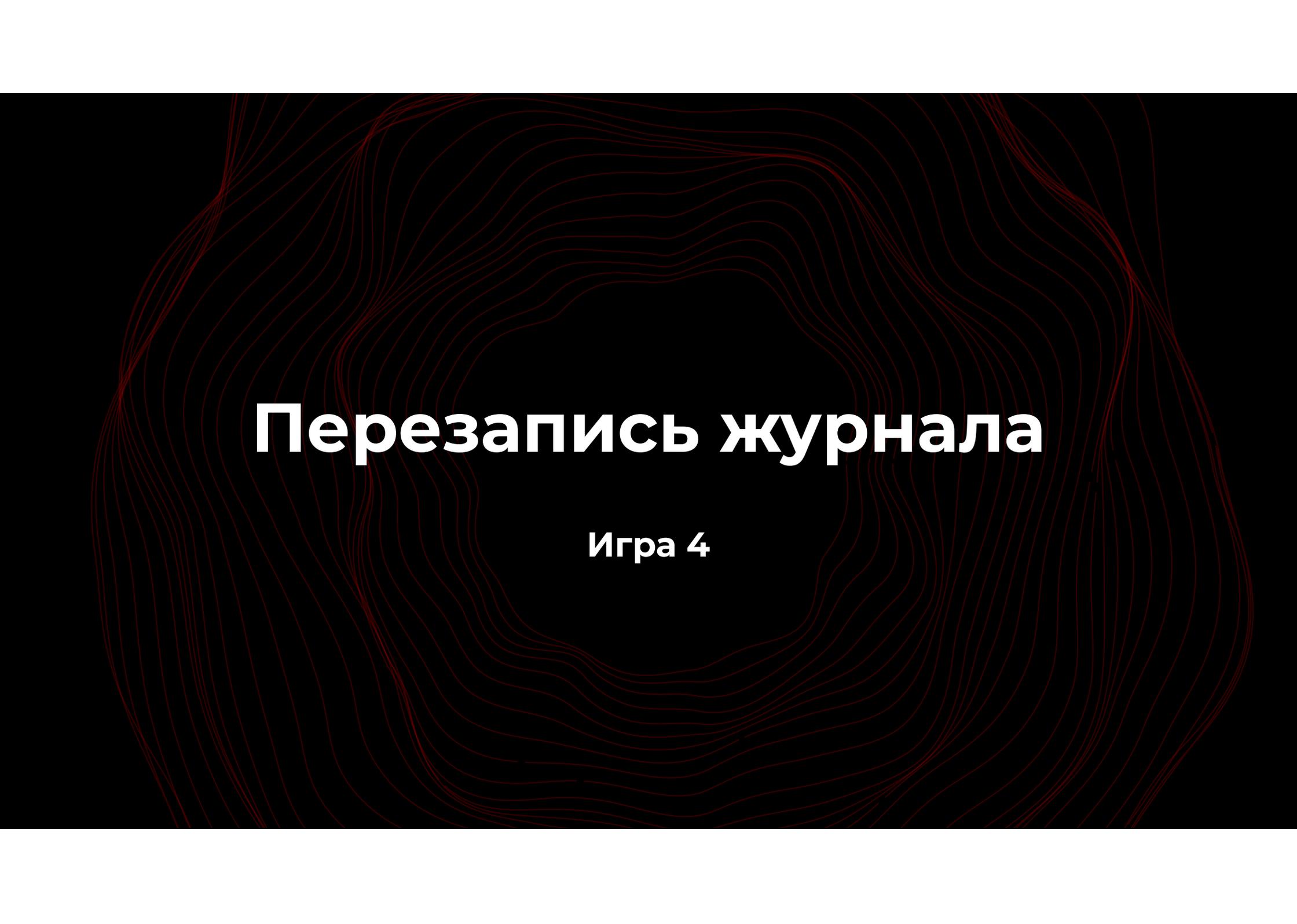
index:	1	2	3	4										
term:	t1	<i>t1</i>	<i>t1</i>	...										
value:	H	<i>i</i>	g	h										
applied:														

RequestVoteRPC

term:	t2
candidateId:	i8
lastLogIndex:	4
lastLogTerm:	t1

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok		X	X		✓	✓	✓	✓



Перезапись журнала

Игра 4

it все еще offline

Потерпите, через 4 слайда вернетесь.

i4, начинайте выборы



Raft State

raft_id
<i>i4</i>

Terms

term:	t1	t2	<i>t3</i>	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>	-	<i>i4</i>											

Raft Log

index:	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>										
term:	<i>t1</i>	<i>t1</i>	<i>t1</i>	...										
value:	<i>H</i>	<i>i</i>	<i>g</i>	<i>h</i>										
applied:														

RequestVoteRPC

term:	<i>t3</i>
candidateId:	<i>i4</i>
lastLogIndex:	<i>4</i>
lastLogTerm:	<i>t1</i>

Response

	i1	i2	i3	<i>i4</i>	i5	i6	i7	i8
ok				<input checked="" type="checkbox"/>				



Вжух, и i4 — лидер

Raft State

raft_id
<i>i4</i>

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	<i>i1</i>	-	<i>i4</i>											

Raft Log

index:	1	2	3	4										
term:	t1	t1	t1	...										
value:	H	i	g	h										
applied:														

RequestVoteRPC

term:	<i>t3</i>
candidateId:	<i>i4</i>
lastLogIndex:	<i>4</i>
lastLogTerm:	<i>t1</i>

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok		X	X	✓	✓	✓	✓	✓



i4, заполняйте raft-журнал

Raft State

raft_id
i4

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	i1	-	i4											

Raft Log

index:	1	2	3	4	5	6	7	8	9	10	11	12		
term:	t1	t1	t1	...	t3	t3	t3		
value:	H	i	g	h	␣	v	o	l	t	a	g	e		
applied:														

AppendEntriesRPC

term:	t3
leaderId:	i4
index:	5 6 7 8
value:	␣ v o l
leaderCommit:	0

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok				✓				



i5-i8, обработайте запрос

AppendEntriesRPC

term:	t3			
leaderId:	i4			
index:	5	6	7	8
value:	␣	v	o	l
leaderCommit:	0			

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok				✓	✓			

Raft State

raft_id
i5

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	i1	-	i4											

Raft Log

index:	1	2	3	4	5	6	7	8						
term:	t1	t1	t1	...	t3	t3	t3	...						
value:	H	i	g	h	␣	v	o	l						
applied:														





іІ, возвращайтесь онлайн

AppendEntriesRPC

term:	t3
leaderId:	i4
index:	5 6 7 8
value:	␣ v o l
leaderCommit:	0

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓			✓	✓	✓	✓	✓

Raft State

raft_id	i1													
---------	----	--	--	--	--	--	--	--	--	--	--	--	--	--

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	i1	-	-											

Raft Log

index:	1	2	3	4	5	6	7	8	5	6	7	8		
term:	t1	t1	t1	t3	t3	t3	...		
value:	H	i	g	h	L	o	a	d	␣	v	o	l		
applied:	✓	✓	✓	✓										

i4 получает ответ



Raft State

raft_id
i4

Terms

term:	t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12	t13	t14
vote:	i1	-	i4											

Raft Log

index:	1	2	3	4	5	6	7	8	9	10	11	12		
term:	t1	t1	t1	...	t3	t3	t3		
value:	H	i	g	h	␣	v	o	l	t	a	g	e		
applied:	✓	✓	✓	✓	✓	✓	✓	✓						

AppendEntriesRPC

term:	t3			
leaderId:	i4			
index:	5	6	7	8
value:	␣	v	o	l
leaderCommit:	0			

Response

	i1	i2	i3	i4	i5	i6	i7	i8
ok	✓	✓	✓	✓	✓	✓	✓	✓

Факультатив

- Динамическое изменение топологии
- Pre-vote
- Снапшоты
- Quorum read

Материалы

Слайды:

- Online: <https://rosik.github.io/2023-highload>
- PDF: [slides.pdf](#)
- Раздатка: [form.pdf](#)

Picodata: <https://picodata.io/>, [@picodataru](#)

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

Обратная связь:



<https://conf.ontico.ru/online/hl2023/lectures/5071966>

