

# Rekorde ...werden immer seltener

- Wieviele Rekorde gibt es eigentlich?
  - z.B. heissester August in einem Jahrhundert...
  - Annahme: keine Korrelation (d.h. zufällige Permutation vorausgesetzt)

- Betrachte Bitfolge:  $i$ -te Stelle 1 gdw.  $i$ -ter Wert ein Rekord

Bsp: 1001011000000100000000001000000000000000  
 (Rekorde traten im 1., 4., 6., 7., 14.,... Jahr auf)

- Wieviele 1en bei gegebener Länge  $n$ ?

Antwort: - damit  $i$ -tes Bit auf 1 steht, muss der  $i$ -te Zahlenwert der Folge  $\geq$  alle ersten  $i$  Zahlenwerte sein  
 - Wahrscheinlichkeit dafür ist  $1/i$  (stimmt's? wieso?)

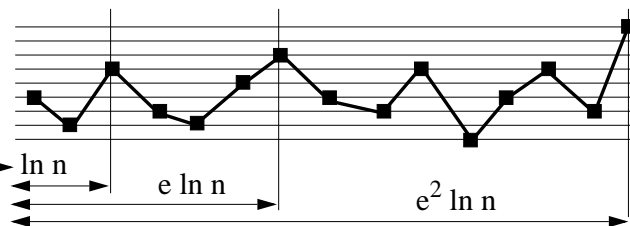
--> Mittlere Anzahl von 1en =  $1 + 1/2 + 1/3 + \dots$   
 - darf man wirklich einfach so aufaddieren?

--> Mittlere Anzahl von Rekorden ist  $H_n \approx \ln n$

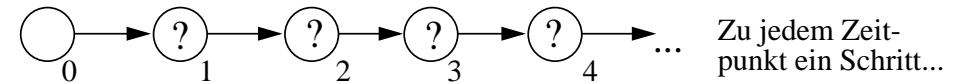
- Bsp: es gibt ca. 5 heisseste Auguste im Jahrhundert  
 ... und wieviele in einem Jahrtausend?

==> Für einen Rekord mehr (im Mittel) muss das Zeitintervall ca.  $e = 2.71828$  mal länger sein

Das haben wir noch immer nicht wirklich bewiesen!



# Verteilung der Lebensdauer



- Wahrscheinlichkeit, *genau* bis zur Position  $i$  zu gelangen (und dort besiegt zu werden)?
- $p(\text{Lebenszeit} = i) = p(\text{Lebenszeit} \geq i) - p(\text{Lebenszeit} > i)$ 
  - Vgl.: 18% Studis  $\geq 16$ . Semester; 13% Studis  $\geq 17$ . Semester  
 ==> 5% *im* 16. Semester > 16. Semester

Wahrscheinlichkeit für "grösster unter den ersten  $i$ "

- Also:  $p(\text{Lebenszeit} = i) = \frac{1}{i} - \frac{1}{i+1} = \frac{1}{i(i+1)}$  (für  $i < n$ )
- Aber:  $p(\text{Lebenszeit} = n) = \frac{1}{n}$  ("erfolgreicher Durchmarsch")

Beispiel  $n=4$ :

gewichtete Summe =  $25/12$  -->  $25/3 = 8.333\dots$   
 Nachrichten im Mittel (vgl. früheres Ergebnis!)

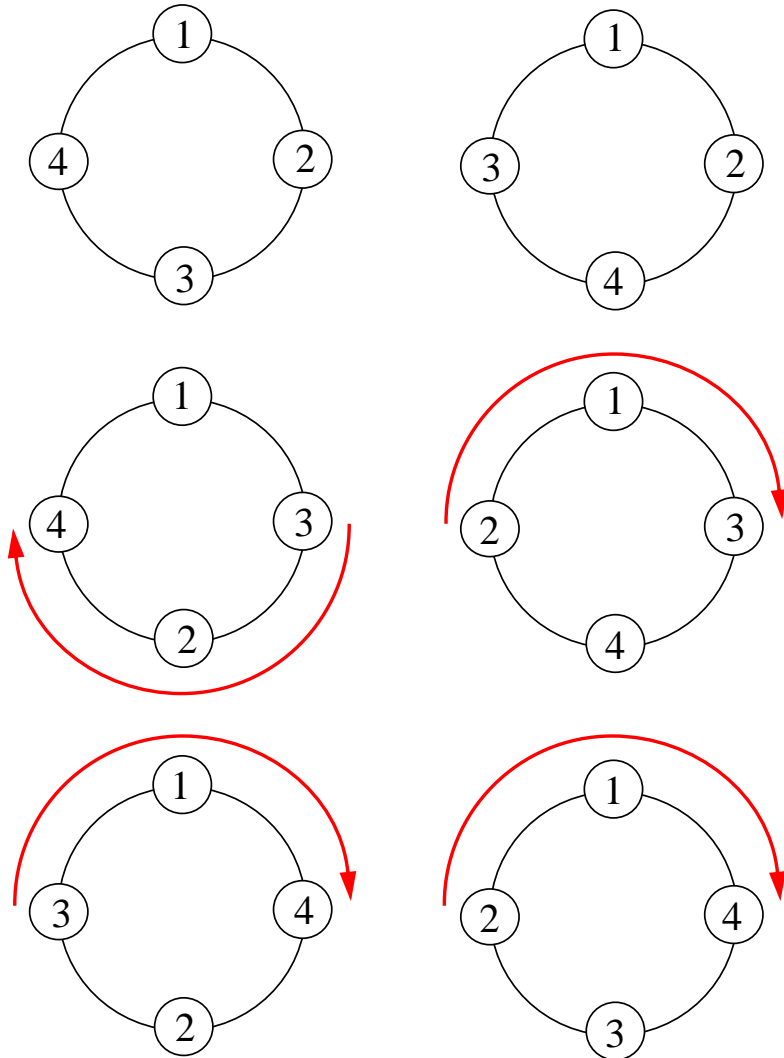
1	$\frac{1}{1 \times 2} = \frac{1}{2} = 50\%$
2	$\frac{1}{2 \times 3} = \frac{1}{6} = 16.7\%$
3	$\frac{1}{3 \times 4} = \frac{1}{12} = 8.3\%$
4	$\frac{1}{4} = 25\%$

vgl. dazu nächstes Bild

100%

Bem.: Es gilt  
 $\frac{1}{n} + \sum_{i < n} \frac{1}{i(i+1)} = 1$   
 (Induktionsbeweis als einfache Übung)

## Lebensdauer 2 bei n=4



6 Bilder

jeweils 4 Starter

In 4 von  $6 \times 4 = 24$  Situationen, also  $1/6$  aller Fälle (= 16.7%), wird für einen Initiator die Lebensdauer 2 exakt erreicht

## Induktionsbeweis der Summenformel

Behauptung: 
$$\sum_{i=1}^{n-1} \frac{1}{i(i+1)} = 1 - \frac{1}{n}$$

Beweis durch vollständige Induktion:

Induktionsanfang:  $n = 1 \rightarrow 0 = 0$  ✓

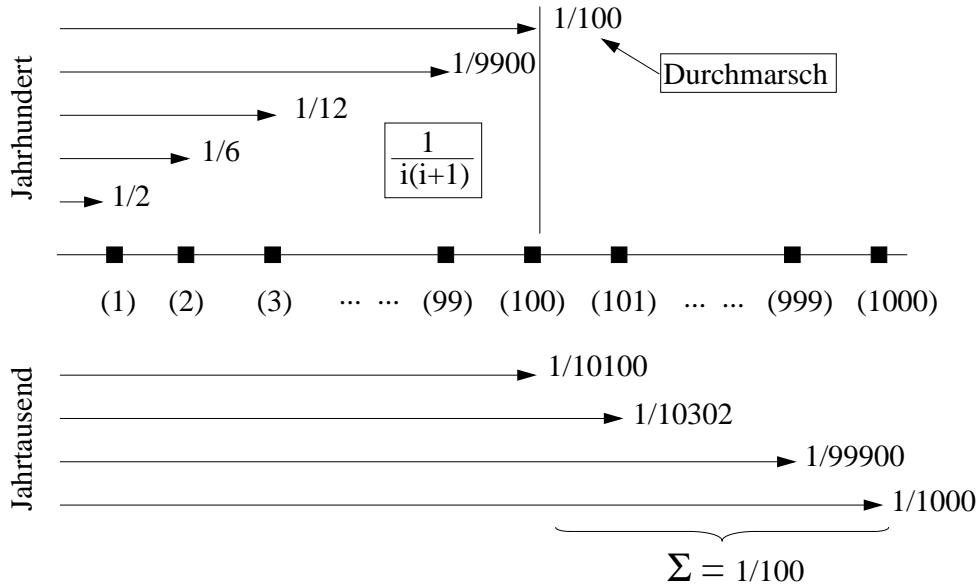
$n = 2 \rightarrow \frac{1}{1 \cdot 2} = 1 - \frac{1}{2}$  ✓

Induktionsschritt:

$$\begin{aligned} \sum_{i=1}^n \frac{1}{i(i+1)} &= \overbrace{1 - \frac{1}{n}}^{\text{Ind. annahme}} + \frac{1}{n(n+1)} \\ &= 1 + \frac{-(n+1) + 1}{n(n+1)} \\ &= 1 + \frac{-n}{n(n+1)} \\ &= 1 - \frac{1}{n+1} \quad \checkmark \end{aligned}$$

# Rekord im Jahr i eines Jahrhunderts?

- Oder: Wahrscheinlichkeit, bei Position (i) "unterzugehen":



- "Durchmarsch" ist wahrscheinlicher, als an der Position davor "unterzugehen"

- Nur 50% Chance, über den ersten hinwegzukommen

-  $p(\text{kein echter Rekord in einem Jahrhundert}) = 1/100$ ,

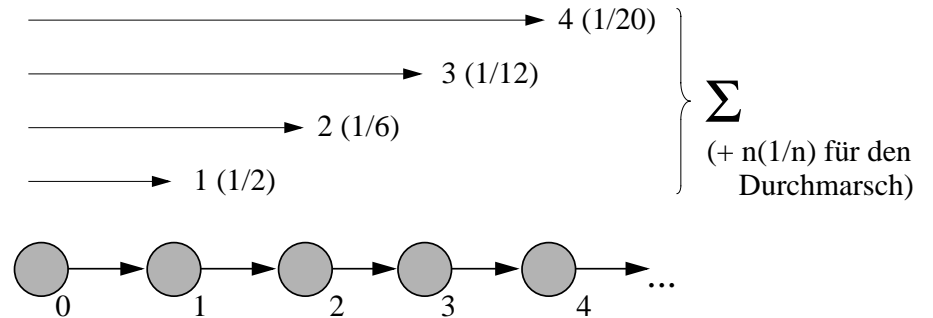
-  $p(\dots \text{ in einem Jahrtausend}) = 1/1000$

Rekord würde jenseits liegen:  
Geht bei Jahrhundert "nur" mit einem Gewicht von 100 ein, statt mehrere hundert bei Jahrtausend

# Nun also: Wartezeit bis zum ersten Rekord

- Gesucht: Die *mittlere* Lebensdauer

- Dazu: Wahrscheinlichkeit mit Weglänge gewichten und alle Einzelfälle aufsummieren



- Daraus folgt allgemein für den Erwartungswert:

$$n \frac{1}{n} + 1 \frac{1}{1 \times 2} + 2 \frac{1}{2 \times 3} + 3 \frac{1}{3 \times 4} + \dots + (n-1) \frac{1}{(n-1)n} = \underline{\underline{H_n}}$$

Also: Eine von einem "x-beliebigen" Prozess initiierte Nachricht kommt im Mittel  $H_n$  weit

- darf man das nun mit n multiplizieren, da dies jeder Prozess tut?
- ja; aber generell Vorsicht bei solchen Überlegungen!

Damit haben wir also  $nH_n$  als mittl. Nachrichtenkomplexität!

## Qualität des Chang/Roberts-Algorithmus

- *Satz* (o. Bew.):

Für unidirektionale Ringe der Grösse  $n$  und  $n$  Initiatoren ist  $n H_n$  die optimale mittlere Nachrichtenkomplexität beim *Election-Problem*

- mit anderen Worten: es gibt keinen anderen Algorithmus, der "besser" ist!

--> Chang/Roberts-Algorithmus ist "average case optimal"!

- *Satz* (Rotem et al., o. Bew.):

Für die Nachrichtenzahl  $m$  des Chang/Roberts-Algorithmus gilt für alle  $\epsilon$ :  $\lim_{n \rightarrow \infty} \text{prob}(m < (1 + \epsilon) n H_n) = 1$

*Interpretation:* Bei grösseren Ringen fast nie mehr als  $n H_n$  Nachrichten...

## Tschebyscheff- / Chernoff-Ungleichungen

- Die Nachrichtenzahl des Chang/Roberts-Algorithmus lässt sich auch genauer abschätzen:

- aus der *Tschebyscheff-Ungleichung* (engl.: "Chebyshev") folgt:

$$\text{prob}(m \geq (1 + \epsilon) n H_n) \leq \frac{\pi^2/6}{\epsilon^2 H_n^2} \quad (\text{vgl. dazu Lehrbücher zur Wahrscheinlichkeitstheorie})$$

- noch besser lässt sich dies mit der *Chernoff-Ungleichung* abschätzen:

$$\text{prob}(m \geq (1 + \epsilon) n H_n) \leq \exp(-\epsilon^2 n H_n / 3)$$

- Einige Beispiele (Wahrscheinlichkeit, dass die Nachrichtenzahl unter dem angegebenen Wert liegt):

	<i>Tschebyscheff</i>	<i>Chernoff</i>
$\epsilon=1, n=10$	0.31	0.00046
$\epsilon=1, n=100$	0.078	$10^{-67}$
$\epsilon=0.5, n=100$	0.31	$10^{-17}$
$\epsilon=0.1, n=100$	--	0.2

- die Tschebyscheff-Ungleichung ist also recht konservativ!

- auch die Chernoff-Ungleichung ist nur eine Abschätzung; die Nachrichtenzahl liegt also fast immer sehr nahe am Mittelwert!

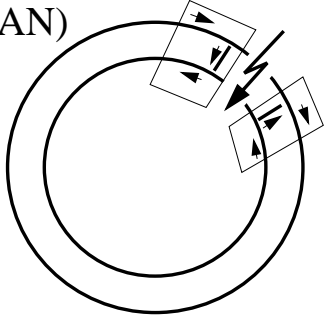
Zur Chernoff-Ungleichung vergleiche man z.B.:

- C. Lavault: Evaluation des algorithmes distribués. Hermes, Paris 1995
- T. Hagerup, C. Rüb: A Guided Tour of Chernoff Bounds. Inform. Proc. Lett. 33, 305-308, 1990

# Fehlertoleranz?

- Was geschieht bei Fehlern? ← z.B. beim Election-Algorithmus
  - was für Fehler können auftreten?
  - welche Fehlerarten betrachtet man? --> *Fehlermodell*
- Es werden typischerweise bei verteilten Systemen verschiedene *Fehlerarten* betrachtet (die unterschiedlich schwer zu behandeln sind), z.B.:
  - *crash* oder *failstop*: Prozess stoppt und bleibt gestoppt
  - *omission* (oder *lossy channel*): Nachricht geht unterwegs verloren
  - *link failure*: Kommunikationsverbindung geht (dauerhaft) kaputt
  - *byzantine*: Prozess verhält sich beliebig falsch (generiert illegale Nachrichten,...)
  - *timing* oder *performance*: bei synchronen bzw. Realzeitsystemen
- *Fehlertolerante* Algorithmen sollen robust bezüglich des betrachteten Fehlermodells sein
- Ob man überhaupt Fehler in Betracht zieht und welches Fehlermodell adäquat ist, ist ein *pragmatischer* Aspekt
  - abhängig von den Anforderungen, der Systemumgebung, Einsatzgebiet...
- Problem der *Fehlerentdeckung*: wie unterscheidet man:
  - Ausfall eines benachbarten Prozesses,
  - Ausfall der Kommunikationsverbindung,
  - extreme Nachrichtenverzögerung (länger als der timeout)?
- Fehlerentdeckung in der Praxis i.a. über *timeouts*
  - z.B. kein Acknowledge innerhalb eines bestimmten Zeitintervalls
  - dann *vermutet* man einen Fehler bzw. *verdächtigt* einen Prozess, ausgefallen zu sein (damit man sich da nicht täuscht, ist es manchmal besser, einen so verdächtigen Prozess nochmals explizit auszuschalten!)

# Fehlerverhalten des Election-Algorithmus?

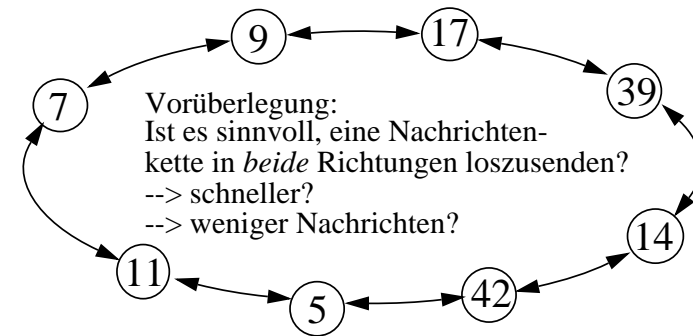
- Bei Ausfall einer Kommunikationsverbindung im Ring kann der Algorithmus nicht mehr funktionieren
  - In der Praxis (z.B. Token-Ring-LAN) verwendet man daher meistens "Doppelringe"
    - Nachbarstationen schliessen den Ring kurz
    - zweite Unterbrechung wird so allerdings nicht toleriert!
- 
- Konsequenz von Prozessausfällen?
    - in der Praxis hat das Ringinterface einer Station ein "Bypass-Relais" (damit wird die ausgefallene Station einfach kurzgeschlossen)
  - Aber: Was geschieht bei folgendem Szenario?
    - Prozess mit der höchsten Identität sendet eine Nachricht aus
    - crasht (und wird kurzgeschlossen), bevor er "seine" Antwort vom Ring nehmen kann
    - diese kreist nun ständig weiter --> Algorithmus terminiert nicht!
  - Denküben:
    - wie kann man obigen Fall verhindern (den Algorithmus also tolerant gegenüber diesem Fall machen)?
    - was geschieht, wenn ein anderer Prozess entsprechend ausfällt?
    - wie macht man den Algorithmus fehlertolerant gegenüber sporadisch verlorenen Nachrichten?

# Praxisrelevanz

- Für die Praxis sind diverse Aspekte von Algorithmen relevant
  - z.B. Fehlertoleranz im o.g. Sinne
  - oder: Crash eines Prozesses, der mitten im Senden ist (was geschieht mit Nachrichtenbruchstücken?)
  - oder: Erkennen verfälschter Nachrichten
- Denkübung: Kann beim Election-Verfahren der Fall behandelt werden, dass zwei verschiedene Stationen (versehentlich) die gleiche Identität haben?
- Der Election-Algorithmus muss wiederholt ausführbar sein
  - was geschieht, wenn in einer Ausführung des Algorithmus alte (sehr langsame) Nachrichten einer früheren Ausführung "dazwischenfunken"?
  - wie lässt sich dieses Problem behandeln?

==> Einfache algorithmische Ideen werden in der Praxis oft in komplexe Algorithmen eingebettet, um den ganzen pragmatischen Ansprüchen gerecht zu werden!

# Probabilistisches Election-Verfahren für bidirektionale Ringe



- Message-extinction-Prinzip in naheliegender Weise verallgemeinern:

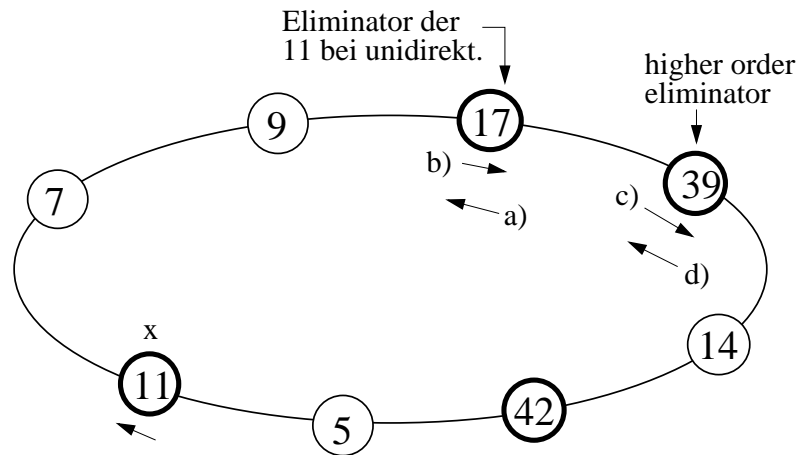
```

Ip: {M = 0}
M := p;
Wähle mit Wahrscheinlichkeit 1/2 eine Richtung;
send <M> to Nachbar in diese Richtung ;

Rp: {Eine Nachricht <j> ist eingetroffen}
if M < j then
    M := j;
    send <M> to ... /* an anderen Nachbarn */
fi
/* weitersenden */
if j = p then "I am the master" fi
    
```

- ist Wahrscheinlichkeit 1/2 eine gute Wahl?
- geht es nicht besser deterministisch als probabilistisch?
- wie hoch ist die mittlere Nachrichtenkomplexität?

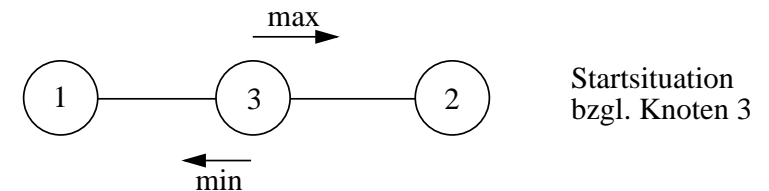
# Mittlere Nachrichtenkomplexität



- Annahme: Jede Einzelnachricht braucht gleich lang
- In der Hälfte aller Fälle kommt der Eliminator der Nachricht "x" auf halbem Weg entgegen
  - z.B. Fall a) bei den Knoten 11 und 17
  - > sollte 1/4 aller Nachrichten (gegenüber unidirekt. Fall) sparen (wieso?)
- Aber: höchste läuft immer ganz durch  
(Jedoch: spielt für  $n \rightarrow \infty$  eine "asymptotisch geringe" Rolle)
- Asymptotische mittlere Nachrichtenkomplexität ist *geringer* als  $0.75 n \ln n$  (Grund: "Higher order eliminators")
  - z.B. 39 verkürzt den Weg der 11 "etwas" im Fall b), d)
  - 42 als higher order eliminator würde hier aber nichts nützen!

# Deterministische bidirektionale Verfahren

- (1) "Gerader" Prozess startet im Uhrzeigersinn, andere gegen der Uhrzeigersinn
  - "common sense of orientation" vorhanden?
  - aber: wenn Identifikationen keine ganze Zahlen (sondern z.B. rationale)?
- (2) Starten in Richtung des kleineren Nachbarn ("min")
  - kennt man Identität der Nachbarn? (wenn nicht, was dann?)
  - wieso nicht in Richtung des *grösseren* Nachbarn ("max")?
  - was bringen die deterministischen Verfahren "min", "max" im Vergleich zum probabilistischen Verfahren?

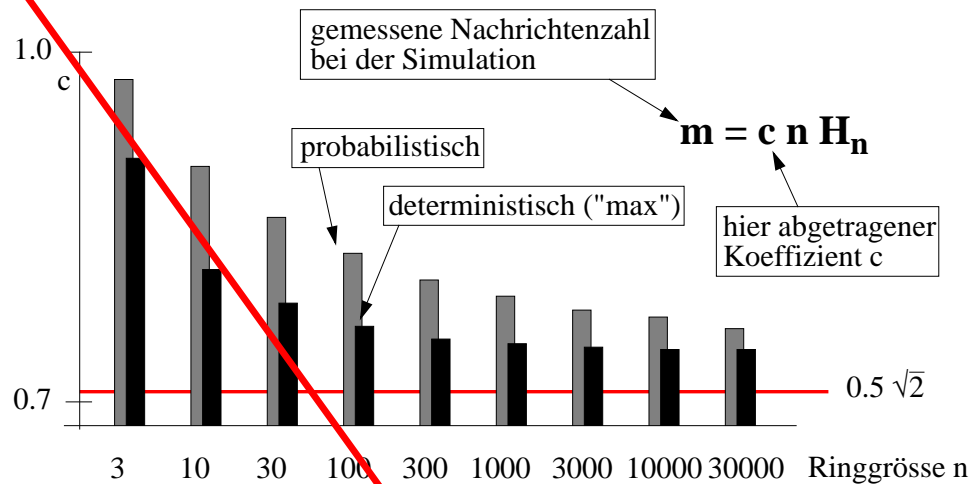


Resultat von Lavault (Beweis schwierig!):  
Asymptot. Nachrichtenkomplexität  $n \rightarrow \infty$  ist

$$0.5 \sqrt{2} n \ln n \approx 0.7071 n \ln n$$

Als *untere Schranke* für das bidir. Election-Problem kennt man  $0.5 n \ln n$  ( $\rightarrow$  "Lücke")

# Gemessene Nachrichtenkomplexität



## - Simulationsergebnisse:

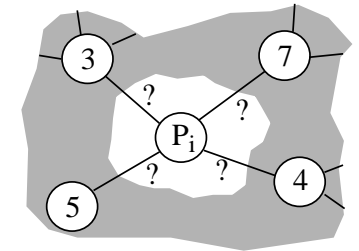
- Es stellt sich heraus, dass "min" und "max" bzgl. der mittleren Nachrichtenkomplexität etwa identisch sind (!), dass diese Varianten jedoch etwas besser sind als die probabilistische Version (und diese, wie gezeigt, besser als die unidirektionale)
- Simulationen zeigen auch, dass die asymptotische Nachrichtenkomplexität des probabilistischen Verfahrens lediglich ein theoretisches Ergebnis ist und der Faktor 0.7071... sehr langsam (mit steigendem n) approximiert wird
- Ferner zeigen die Simulationen, dass die *Abweichungen vom Mittelwert*  $n H_n$  bzgl. der Nachrichtenkomplexität i.a. nur sehr gering sind; 100000 Simulationen bei einer Ringgröße von 20 lieferten z.B. stets Nachrichtenzahlen unter  $2 n H_n$
- Beachte bei *Simulationsexperimenten*: sehr viele Einzelexperimente (Varianz, statistisch relevante Ergebnisse) sowie guter Zufallszahlengenerator notwendig

# Nachbarschaftswissen

- Zweck: Lerne Identitäten Deiner Nachbarn kennen (Idee: "Ich heiße i, wie heißt Du?")

- Hier nachrichtengesteuert für  $P_i$ : (vgl. verteiltes Approximationsschema)

{ Nachricht  $\langle j \rangle$  kommt an  
**if not vorgestellt then**  
**send  $\langle i \rangle$  to all neighbors;**  
**fi;**  
 Nachbarn := Nachbarn  $\cup \{j\}$ ;  
 vorgestellt := **true**;



Netz zusammenhängend  
 mit bidirektionalen Kanten

mehrere?  
 Einer muss "spontan" mit  
 der Vorstellung beginnen

- *global terminiert*, wenn
  - vorgestellt = true bei allen Prozessen und keine Nachricht unterwegs
  - *oder*: alle kennen alle
- *lokal terminiert* (genügt meistens!), wenn zu jeder inzidenten Kante die zugehörigen Nachbarknoten bekannt sind
- $2e$  Nachrichten (bei Ring also  $2n$ ); Zeitkomplexität?

## Verwendung als:

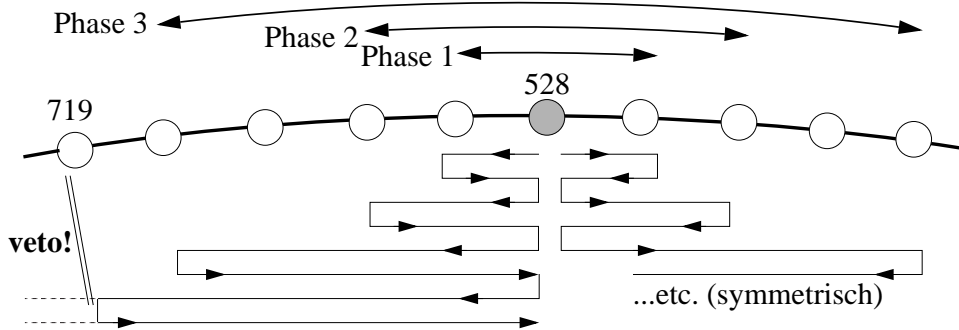
- vorgeschalteter eigener Algorithmus
- "piggybacking" der Vorstellungsnachrichten mit eigentlichen (jew. "ersten") Anwendungsnachrichten
- lohnt sich das bei bidirektionaler Election, um so den als "min" oder "max" bezeichneten Algorithmus einsetzen zu können?



# Hirschberg / Sinclair-Election-Algorithmus

- Idee: Jeder Knoten versucht, sukzessive Gebiete der Grösse  $2^i$  ( $i=1, \dots$ ) zu erobern

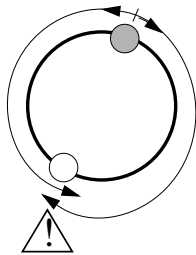
bidirektionaler Ring!



- Ein unterwegs angetroffener grösserer Knoten legt Veto ein
  - > Initiator über Rückmeldung informieren
  - > Initiator wechselt von *aktiv* nach *passiv*

nur noch Nachrichten weiterleiten ("relay")

Gewinnsituation:

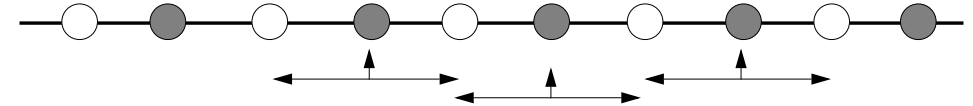


- Nachricht läuft in bereits selbst erobertes Gebiet
- oder: Nachricht trifft bei Initiator selbst wieder ein
- es bleibt genau ein Gewinner! (wieso?)

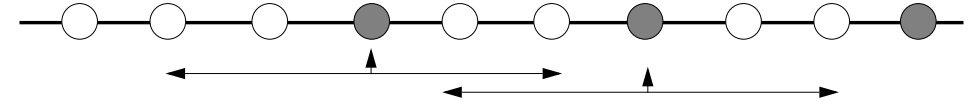
# Komplexitätsanalyse

- Ein Prozess kann nur dann eine Kette der Länge  $2^i$  starten, wenn er im Abstand  $2^{i-1}$  in beiden Richtungen überlebt hat
  - Dichte überlebender Prozesse nimmt also exponentiell ab
- Innerhalb eines Bereiches von  $1 + 2^{i-1}$  benachbarter Prozesse kann also höchstens einer eine Kette der Länge  $2^i$  starten

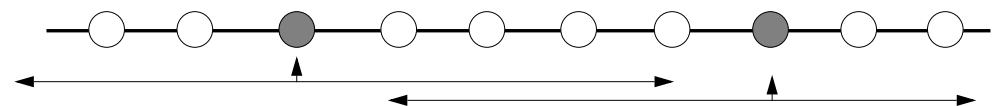
- Jeweils 1 dazwischenliegender Prozess nach Phase 1:



- Jeweils 2 dazwischenliegende Prozesse nach Phase 2:



- Jeweils 4 dazwischenliegende Prozesse nach Phase 3:



- $n/2$  Prozesse können Ketten der Länge 2 initiieren
- $n/3$  Prozesse können Ketten der Länge 4 initiieren
- $n/5$  Prozesse können Ketten der Länge 8 initiieren
- ...
- $n/(1+2^{i-1})$  Prozesse können Ketten der Länge  $2^i$  initiieren

## Maximal $8 n \log_2 n$ Nachrichten

- Also: höchstens  $n/(1+2^{i-1})$  Prozesse initiieren eine Nachrichtenkette der Länge  $2^i$  in Phase  $i$
- Bei jeder solchen Kette wird jede Kante max. 4 Mal durchlaufen
- In Phase  $i$  gibt es also höchstens  $4 \times 2^i \times n / (1+2^{i-1}) < 8n$  Nachrichten
- Es gibt höchstens  $1 + \lceil \log_2 n \rceil$  Phasen

$\implies$  ca.  $8 n \log_2 n \approx 5.55 n \ln n$  Nachrichten *maximal*  
(Worst-case-Komplexität!)

- 
- Aber: Wie hoch ist die *mittlere* Nachrichtenkomplexität?
  - *Zeitkomplexität*:  $2 + 4 + 8 + 16 + \dots + 2^i < 2^{i+1} \approx 4n$

---

Vergleiche dies alles mit dem Chang / Roberts - Algorithmus

- welcher Algorithmus ist "in der Praxis" besser?

## Synchrone $\leftrightarrow$ asynchrone Phasen

- Die Phasen der einzelnen Initiatoren müssen nicht unbedingt synchron laufen!
- Damit der Algorithmus dann noch gut funktioniert, vereinbare folgendes:
  - anstelle von Knotenidentitäten betrachtet man das Paar (Phasennummer, Knotenidentität)
  - diese Paare werden lexikographisch geordnet (eine höhere Phasennummer hat also Priorität!)
- Konsequenz: Ein "schneller" Initiator gewinnt gegenüber einem "langsamen" Initiator mit höherer Identität
- Es gewinnt also zwar ein eindeutiger Knoten die Election, das muss aber nicht derjenige mit der grössten Identität sein!

- 
- Unterscheide also:

- *leader election problem*
- *maximum finding problem*



- eine Lösung des maximum finding problems ist immer auch eine Lösung des leader election problems (sofern die Knoten eindeutig nummeriert sind)
- Umkehrung?