

# Datastructuren en Algoritmen III, 2016

Gunnar Brinkmann

19 september 2016

## Inhoudsopgave

<b>1</b>	<b><u>Introductie</u></b>	<b>2</b>
<b>2</b>	<b><u>Metaheuristieken</u></b>	<b>3</b>
2.1	<u>Local search</u> . . . . .	5
2.1.1	<u>Het belang van de buurfunctie</u> . . . . .	9
2.2	<u>Guided local search</u> . . . . .	12
2.3	<u>Variable neighbourhood metaheuristieken</u> . . . . .	17
2.4	<u>Simulated annealing</u> . . . . .	21
2.5	<u>Genetische algoritmen</u> . . . . .	26
2.6	<u>Afsluitende opmerkingen</u> . . . . .	33
<b>3</b>	<b><u>Algoritmen voor slim gebruik van het geheugen</u></b>	<b>34</b>
3.1	<u>Hashing</u> . . . . .	35
3.1.1	<u>Extendible hashing</u> . . . . .	36
3.1.2	<u>Linear hashing</u> . . . . .	42
3.1.3	<u>Bloom filters</u> . . . . .	48
3.2	<u>Sorteren van grote hoeveelheden data – extern sorteren</u> . . . . .	53
3.3	<u>Grote zoekbomen</u> . . . . .	56
3.3.1	<u>B-trees</u> . . . . .	57
3.3.2	<u>B+-trees</u> . . . . .	63
<b>4</b>	<b><u>Algoritmen voor strings</u></b>	<b>69</b>
4.1	<u>Exact string matching</u> . . . . .	69
4.2	<u>Benaderend string matching</u> . . . . .	95
<b>5</b>	<b><u>Compressiealgoritmen</u></b>	<b>106</b>
5.1	<u>Huffman codering</u> . . . . .	110
5.2	<u>LZ77</u> . . . . .	116

5.3	<u>LZW</u> . . . . .	121
5.4	<u>De Burrows-Wheeler transformatie.</u> . . . . .	127
<b>6</b>	<b><u>Parallele algoritmen</u></b>	<b>136</b>
6.1	<u>Branch and bound met verdeelde processoren</u> . . . . .	139
6.2	<u>Op weg met MPI</u> . . . . .	145
6.2.1	MPI opstarten . . . . .	146
6.2.2	Berichten sturen en ontvangen . . . . .	149
6.2.3	Berichten sturen en ontvangen – zonder de voortgang te blokkeren . . . . .	156
6.2.4	Verder met MPI? . . . . .	159
6.3	<u>De modellen voor parallel computing</u> . . . . .	159
6.4	<u>Semigroep algoritmen</u> . . . . .	166
6.5	<u>Eén parallel grafen algoritme</u> . . . . .	168
6.6	<u>Pipelining</u> . . . . .	178

# 1 Introductie

De bedoeling van deze tekst is **niet** in plaats van de les te worden gebruikt maar alleen om te helpen de eigen nota's misschien een beetje beter te verstaan als er iets niet echt duidelijk is. Alleen in de les kan je onmiddellijk vragen stellen als je iets niet verstaat...

## 2 Metaheuristieken

Veel problemen kunnen niet optimaal opgelost worden – bv. NP-complete problemen voor een grote input. Jammer genoeg geldt dat ook voor veel problemen die in de praktijk belangrijk zijn! Aan de universiteit kan je dan misschien zeggen dat je liever een ander probleem bestudeert, maar in *het echte leven* is dat niet altijd mogelijk. Als je bv. in een bedrijf werkt en een optimale manier moet vinden om een machine te besturen, dan kan je soms aantonen dat het probleem heel moeilijk is en dat je weinig kans hebt een optimale oplossing te vinden. Jouw baas dan voorstellen gewoon een ander probleem op te lossen dat gemakkelijker is, is misschien niet echt een goed idee. . .

Wij moeten in zulke gevallen dus gewoon proberen het beste te doen wat we kunnen – ook al vinden wij geen optimale oplossing. Wij moeten dus heuristieken toepassen om ten minste aanvaardbare oplossingen te vinden. De taak van metaheuristieken is ons te helpen dergelijke heuristieken te ontwikkelen. Er bestaan ongelofelijk veel verschillende metaheuristieken. In deze les gaan wij een heel kort en oppervlakkig overzicht geven van verschillende – maar zeker niet alle – metaheuristieken. De reden dat wij maar een overzicht geven en niet echt alle details bespreken is dat als jullie eens in een situatie zouden zijn waar jullie zo’n metaheuristiek moeten toepassen – dan moeten jullie er toch heel veel tijd aan besteden om te zien welke trucs voor deze metaheuristiek al gebruikt werden, hoe je hem het best op dit specifieke probleem kan toepassen, de parameters optimaliseren, etc. Dus: Je moet veel meer weten dan men in één enkele les kan vertellen. En als jullie nooit zelf een heuristiek moeten ontwikkelen dan hebben jullie de details natuurlijk niet nodig. Maar in beide gevallen is het nuttig sommige ideeën te kennen en ongeveer te weten wat er bestaat.

Een Metaheuristiek is een heuristiek om heuristieken te ontwikkelen. Dus een manier om voor sterk verschillende problemen heuristische algoritmen te ontwerpen. Maar over welk soort problemen hebben wij het:

In principe heb je altijd een (eindige of oneindige) verzameling  $M$  van punten en een doelfunctie (objective function)  $f : M \rightarrow \mathbb{R}$ . Je zoekt het punt  $x$  waarvoor  $f(x)$  maximaal of minimaal is. In principe moet je er natuurlijk over nadenken, of zo’n punt bestaat, maar in de meeste gevallen (bv. als de verzameling eindig is) is dat geen probleem.

**Voorbeeld 1** *Situatie: Je bent op een vreemde planeet – noemen wij hem eens de aarde – waar er heel veel mist is. Je kan nauwelijks jouw voeten zien, maar je kan wel de hoogte meten van het punt waar je staat. Wat doe je om het hoogste punt van deze planeet te vinden?*

*Of met andere woorden:*

*M is de verzameling van alle punten op de aarde en  $f(x)$  is de hoogte van punt  $x$ . Het doel is dus de top van de hoogste berg van de wereld te vinden. Dit is inderdaad helemaal geen slecht voorbeeld omdat het gebruikt kan worden (en later gebruikt zal worden) om sommige problemen en oplossingen beter te verstaan.*

*Bovendien kan je andere problemen als zijnde analoog interpreteren – het is gewoon een andere aarde en een ander begrip van hoogte.*

Welke manieren om maxima van functies te berekenen kennen jullie?

**Voorbeeld 2**  $M = [-100, 100] \subset \mathbb{R}$ ,  $f(x) = -x^2 + 10x$ .

*Hier zou je natuurlijk jouw kennis uit de analyse toepassen en  $f$  afleiden, bepalen waar  $f'(x) = 0$  is en dan kijken of het een maximum of minimum is. Dus  $f'(x) = 10 - 2x$  – dat is 0 voor  $x = 5$  en daar heb je ( $f''(x) = -2$ ) een lokaal en ook globaal maximum. Heel gemakkelijk.*

*Maar hier zien wij meerdere dingen: Wij gebruiken niet alleen maar de verzameling van punten – wij gebruiken ook dat de verzameling een zekere structuur heeft – en dat deze structuur ook iets te maken heeft met de functie (bv. dat voor punten die heel dicht bij elkaar liggen de waarden van de functie ook heel weinig verschillen. Inderdaad gebruiken wij zelfs dat ze continu en afleidbaar is om de technieken te kunnen toepassen!*

*Zo'n probleem is ongeveer het beste geval dat je kan tegenkomen.*

Maar er zijn ook gevallen waar je niet zo veel structuur hebt:

**Voorbeeld 3**  $M = \{(x, y) \in \mathbb{N} \times \mathbb{N} | 0 \leq x, y \leq 1.000.000\}$  en voor elk paar  $(x, y)$  kies je een toevallig getal  $f(x, y) \in \mathbb{R}$ .

*Op het eerste gezicht lijkt dit voorbeeld misschien gemakkelijker omdat je maar eindig veel verschillende waarden hebt – maar hoe vind je hier het maximum ?*

*De methoden uit het vorige voorbeeld kunnen wij hier zeker niet toepassen. Inderdaad is de enige manier om **zeker** te zijn het maximum te hebben gevonden naar **alle** punten te kijken.*

*Ook hier heeft  $M$  een zekere structuur – je kan bv. zeggen dat  $(x, y)$  en  $(x', y')$  buuren zijn als  $|x - x'| + |y - y'| = 1$*

*maar jammer genoeg is er geen samenhang met  $f()$ , dus helpt ze niet met het vinden van het maximum. Als  $f(x, y)$  een hoge waarde heeft, dan is de kans dat voor een buur  $(x', y')$  de functie  $f(x', y')$  ook een grote waarde heeft even groot en even klein als voor paren  $(x', y')$  die veraf liggen.*

*Inderdaad moet je niet naar ingewikkelde strategieën kijken om dit probleem op te lossen. Elke verzameling van  $n$  punten heeft dezelfde kans het maximum*

te bevatten. Als je dus **het** maximum moet hebben, dan moet je toch al alle punten testen en als er (bv. omdat de tijd beperkt is) een bovengrens  $b$  voor het aantal punten is, dat je kan testen, kan je om het even welke  $b$  punten kiezen en het maximum nemen. De kans dat het een goede benadering is, is altijd dezelfde.

Maar zo'n probleem is echt het slechtst mogelijke geval. . .

### Algoritme 1 Toevallig zoeken:

Voorwaarde: gegeven een verzameling  $M$ , een functie  $f : M \rightarrow \mathbb{R}$  en stel dat wij het maximum van  $f()$  zoeken.

Leg een getal  $n$  vast.

- Kies toevallig (en onafhankelijk)  $n$  punten en kies het punt met de grootste waarde van  $f()$ .

Het is vrij duidelijk wat je anders moet doen als je het minimum zoekt. . .

Voor het voorbeeld 2 zou toevallig zoeken zeker niet het optimum vinden. Of het algoritme een aanvaardbare oplossing vindt hangt ervan af hoe vaak je kiest en hoe goed de benadering moet zijn. . . . Het beste algoritme voor dit geval is het duidelijk niet. Voor het voorbeeld 3 gaat het vrij zeker ook geen goede oplossing vinden, maar toch is het algoritme voor dit voorbeeld min of meer optimaal – er bestaan gewoon geen echt betere algoritmen. . . .

**Oefening 1** Hoe denk je dat toevallig zoeken voor het voorbeeld 1 zou presteren? Waarvan hangt het af hoe goed het algoritme presteert ?

Metaheuristieken steunen erop dat – hoewel wij niet in de ideale situatie zitten – het tenminste mogelijk is het zo te modelleren dat de structuur van de verzameling  $M$  waarop wij het optimum zoeken *een beetje vriendelijk* is en een beetje samenhang met de functie vertoont. Inderdaad is de *structuur* van  $M$  normaal geen deel van het probleem en ze moet nog vastgelegd worden als wij zo'n structuur nodig hebben.

Het is het gemakkelijkst om te begrijpen als wij gewoon beginnen naar sommige metaheuristieken te kijken.

## 2.1 Local search

Local search is ongeveer de eenvoudigste – maar zeker niet de slechtste – metaheuristiek. Het is in principe hetzelfde als een *gerandomiseerd gretig algoritme*. Stel dat er voor elk punt in de verzameling  $M$  zekere stappen (of

operaties) bekend zijn om van dit punt naar een ander te gaan. Punten waar je naartoe kan gaan, noemen wij buren en voor de verzameling van alle buren van een punt  $a$  schrijven wij  $N(a)$ . Als je van punt  $a$  naar punt  $b$  in één stap kan komen – dus  $b \in N(a)$ , dan interpreteren wij dat zo, dat op onze kaart  $a$  dicht bij  $b$  ligt.

### Algoritme 2 Local search:

*Gegeven zijn  $M$  en  $f()$  en wij hebben  $N()$  al gekozen, wij weten dus welke punten buren van elkaar zijn. Stel dat wij het maximum zoeken.*

- a.) *begin met een toevallig gekozen punt  $x$ .*
- b.) *kies een toevallige nog niet gekozen buur  $y$  van  $x$  (dus  $y \in N(x)$ ). Als zo'n buur niet bestaat, stop dan.*
  - *als  $f(y) > f(x)$ , kies dan  $y$  als jouw nieuwe positie  $x$  en herhaal b.)*
  - *anders herhaal b.) met jouw oude positie  $x$*

Soms zijn er variaties op dit algoritme nodig:

- Inderdaad moet je in stap b.) niet altijd alle mogelijke stappen testen – het is bv. mogelijk dat er oneindig veel zijn. Je kan ook een bovengrens  $B$  vastleggen en als je na  $B$  stappen nog geen betere positie hebt gevonden, stop je ook.
- Je kan in b.) ook maar  $f(y) \geq f(x)$  in plaats van  $f(y) > f(x)$  eisen als je op een andere manier garandeert dat je altijd gaat stoppen (bv. stop altijd na een zeker aantal stappen of hou de punten bij waar je al was, etc.)

**Oefening 2** *Een gretig algoritme zou bv. op dezelfde manier kunnen beginnen als lokaal zoeken maar voor elk punt niet toevallig een buur kiezen en die aanvaarden als de waarde van  $f()$  beter is, maar altijd de **beste** buur kiezen (als dat mogelijk is – dus bv. als de verzameling van buren eindig is).*

- *Geef een voorbeeldprobleem, een buurfunctie en een startpunt waar lokaal zoeken soms betere resultaten oplevert dan dit gretig algoritme.*
- *Geef een voorbeeldprobleem, een buurfunctie en een startpunt waar lokaal zoeken soms slechtere resultaten oplevert dan dit gretig algoritme.*

Hierbij betekent soms dat dat natuurlijk van de toevallige keuzes van burens kan afhangen en betere oplossing dat het resultaat van de zoektocht een punt met een hogere waarde van  $f()$  is.

**Oefening 3** Beschrijf **precies** een lokaal zoeken algoritme voor het probleem in Voorbeeld 2. Wat kies je als  $N(x)$  voor een punt  $x$ ? Met welke kans kies je een punt  $y \in N(x)$ ? Denk je dat het algoritme goed presteert?

In ons voorbeeld 1 is het duidelijk wat je normaal als burens (dus  $N()$ ) zou kunnen kiezen: bv. alle punten met een afstand van ten hoogste 1 meter van een punt  $x$  zijn burens van  $x$ . De verdeling om toevallig te kiezen zou dan bv. de uniforme verdeling van al deze punten zijn, maar je kan ook een voorkeur geven aan kleine of grote afstanden.

Stel dat je twee buurfuncties  $N_1()$ ,  $N_2$  hebt en voor elk punt  $m \in M$  geldt dat  $N_1(m) \subseteq N_2(m)$ . Als lokaal zoeken dan vanuit een punt  $x$  met buurfunctie  $N_1()$  het optimum kan vinden dan zeker ook met  $N_2()$  – terwijl dat omgekeerd niet noodzakelijk geldt. Het lijkt op het eerste gezicht dus alsof grotere verzamelingen van burens beter zijn... Maar ook al **kan** je het optimum dan ook met  $N_2()$  vinden, is het best mogelijk dat de kans het optimum binnen een beperkt aantal stappen te vinden **zeer veel** kleiner is dan met  $N_1()$ . Om dat goed in te zien is Oefening 6 zeer nuttig. Precies de juiste buurfunctie te definiëren is soms een echte kunst...

**Oefening 4** Waar zou je terecht komen als je bv. in de leszaal zou starten en dit algoritme toepassen om het hoogste punt van de wereld te vinden?

**Oefening 5** Zou het niet beter zijn (het is maar theoretisch...) in Oefening 4 toe te laten dat de stappen t.e.m. 20 000 km lang mogen zijn? Dan kan je zeker zijn altijd kans te maken het optimum te vinden – om het even waar je start.

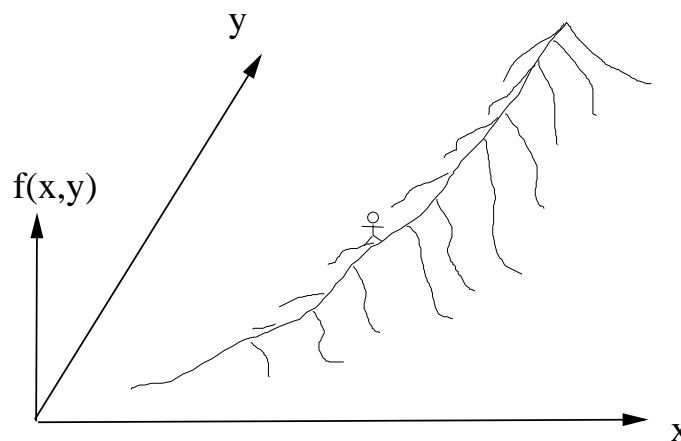
**Oefening 6**  $M = \{(x, y) | x, y \in \mathbb{Z}\}$ , de doelfunctie  $f$  is gedefinieerd door  $f(x, y) = (1 + \frac{1}{1+x^2}) * (1 - (x - y)^2)$ .

Wij definiëren voor elke  $k \in \mathbb{N}$  een buurfunctie:  $N_k(x, y) = \{(v, w) \in M | |v - x| \leq k, |w - y| \leq k\}$ .

- Waar is het maximum van deze functie. Waar is zij positief, waar 0, waar negatief.
- Stel dat je lokaal zoeken in het punt  $x = (5, 5)$  begint.
  - Kan je door middel van lokaal zoeken het maximum vanuit  $x = (5, 5)$  voor elke  $k$  vinden?

- Als je  $N_1$  gebruikt: Wat is in het slechtste geval het aantal stappen dat je nodig hebt, om het maximum te vinden?
- Als je  $N_{20}$  gebruikt: Hoeveel pogingen heb je ongeveer nodig tot de kans een enkel beter punt te hebben gevonden groter is dan  $1/2$ ?
- Stel dat je ver van het optimum bent, maar wel in een punt  $(x, y)$  met  $x = y$  zit. Hoeveel ben je na één stap gemiddeld dichterbij het optimum als je  $N_1()$  gebruikt en hoeveel als je  $N_k()$  voor  $k \rightarrow \infty$  gebruikt? Zou je beter kleine of grote  $k$  gebruiken? Wat is het verschil als je op een vaste afstand van het optimum zit (b.v.  $(x+3, y+3)$  als  $(x, y)$  het optimum is)?
- Stel dat  $M$  niet 2-dimensionaal, maar 3-dimensionaal is en  $f()$  analoog (bv.  $f(x, y, z) = (1 + \frac{1}{1+x^2}) * (1 - ((x-y)^2 + (x-z)^2 + (y-z)^2))$ ). Zij  $N_k(x, y, z) = \{(v, w, t) \in M \mid |v-x| \leq k, |w-y| \leq k, |t-z| \leq k\}$ . Zou je beter grote of kleine  $k$  kiezen voor  $N_k()$ ? Argumenteer eerst zonder te rekenen en test dan of jouw vermoeden klopt.

Hier kan je al zien, hoe belangrijk het is goed te kiezen welke stappen toegelaten zijn: Als je heel lange stappen toelaat, kan je zeker het maximum vinden – maar aan de andere kant heb je zo veel keuze dat het alleen maar toevallig zoeken is en je de structuur van het probleem helemaal niet benut. In de volgende situatie zou je als je voor kleine stappen kiest vrij zeker de weg naar boven vinden. Maar hoe groter de stappen zijn hoe kleiner de kans wordt de top te vinden. Veel kleine stappen zijn hier beter dan weinig grote stappen.



In gevallen waar de structuur van sommige problemen een beetje op de structuur van deze figuur lijkt, is local search een goede oplossing.



### Local search herhalen:

Om meer kans te maken niet alleen een lokaal maar een globaal optimum te vinden, wordt de lokale zoekmethode normaal niet één maar meerdere keren toegepast. Omdat het beginpunt toevallig gekozen wordt, vertrek je normaal elke keer van een ander punt. Dat is inderdaad iets dat voor bijna elke metaheuristiek geldt: Soms is het beter de heuristiek meerdere keren toe te passen en de beste oplossing van de verschillende toepassingen te kiezen dan dure methoden te gebruiken om de resultaten van een enkele toepassing te verbeteren.

#### 2.1.1 Het belang van de buurfunctie

Maar wij moeten eens naar een voorbeeld kijken waar het niet zo duidelijk is op welke manier punten burens zijn en wat de stappen zijn. Hier zullen wij zien hoe buitengewoon belangrijk een goede keuze van de buurfunctie is.

**Voorbeeld 4** *Voor het volgende voorbeeld zou je in de praktijk beter geen metaheuristiek toepassen (Prim, Kruskal), maar het is toch een goed voorbeeld omdat het helpt sommige belangrijke feiten te verstaan:*

*Gegeven een graaf  $G = (V, E)$  en gewichten  $g(e) > 0$  voor de bogen van  $G$ . Wij zoeken een opspannende samenhangende deelgraaf  $B$  van  $G$  zodat de som van de gewichten van de bogen in  $B$  minimaal is.*

*Herinnering: een **opspannende** deelgraaf is een deelgraaf die alle toppen bevat.*

*Dus  $M = \{B \mid B \text{ is opspannende samenhangende deelgraaf van } G\}$  en voor  $B = (V_B, E_B) \in M$  definiëren wij  $g(B) = \sum_{E \in E_B} g(E)$*

*Maar wij hebben nog niet vastgelegd wat de topologie van  $M$  is – dus welke elementen wij als burens willen beschouwen.*

*Eén mogelijkheid:*

**Definitie: (a)** *Voor een opspannende deelgraaf  $B = (V, E)$  definiëren wij  $N_a(B) = \{B' = (V', E') \mid B' \text{ is opspannende deelgraaf en } \#((E - E') \cup (E' - E)) \leq 1\}$ . Twee grafen  $B_1, B_2 \in M$  zijn dus burens, als ze in maar 1 boog verschillen. Dat betekent dat de ene graaf het resultaat van het toevoegen (of verwijderen) van een boog tot (van) de andere is.*

*Nu hebben wij alles om lokaal zoeken toe te passen. Wij beginnen met een toevallig gekozen startgraaf die samenhangend is en passen dan lokaal zoeken toe. Werkt dat goed?*

*Het is inderdaad niet duidelijk hoe je op een toevallige manier een samenhangende deelgraaf moet kiezen – maar dat is een ander probleem.*

*Stel dat  $\bar{B}$  een optimale oplossing is en  $B_0$  de toevallig gekozen startgraaf. Kan het optimum vanuit  $B_0$  worden bereikt? Als  $\bar{B} \subset B_0$ , kan het door het verwijderen van overbodige bogen zeker bereikt worden. Elke keer je een boog verwijdert wordt  $g()$  kleiner. Maar als  $\bar{B}$  ook maar één enkele boog bevat die niet in  $B_0$  zit dan is er ten minste één stap waar een boog toegevoegd moet worden – en dat zou leiden tot een graaf met een grotere waarde van  $g()$  – dus zou in het geval  $\bar{B} \not\subset B_0$  het optimum niet gevonden worden.*

**Definitie: (b)**

$N_b(B) = \{B' = (V', E') \mid B' \text{ is opspannende deelgraaf en } \#(E - E' \cup E' - E) \leq 2\}$ . Twee grafen  $B_1, B_2 \in M$  zijn dus burenen, als ze in maar 1 of 2 bogen verschillen.

*In het geval dat  $B_1, B_2$  bomen zijn, is het gemakkelijk in te zien (maar ook een leuke oefening...) dat de twee bogen  $e_1, e_2$  waarin ze verschillen op een cykel in  $B_1 \cup B_2$  moeten liggen en de ene boog alleen maar in  $B_1$  zit en de andere alleen maar in  $B_2$ .*

*Je kan aantonen (zie oefeningen) dat local search met definitie (b) vanuit elke samenhangende startgraaf  $B_0$  het globale optimum kan vinden.*

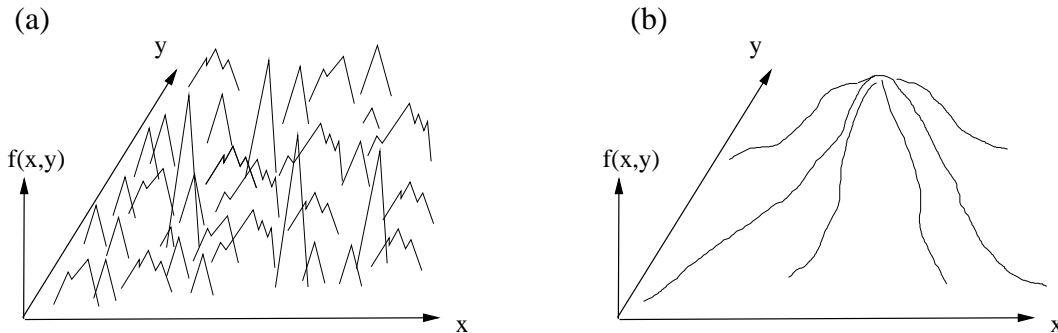
**Opmerking 1** *De verzameling van punten waarvan wij het optimum moeten vinden en de functie zijn deel van het probleem – daarop heb je geen invloed. De definitie van burenen is deel van de methode die je toepast en kan je dus kiezen.*

*Dit voorbeeld toont aan dat het buitengewoon belangrijk is op een goede manier te kiezen wat burenen zijn. Dat geldt niet alleen voor lokaal zoeken maar ook voor andere metaheuristieken die wij zullen zien. Soms is het helemaal niet gemakkelijk uit te vissen wat een goede manier is om  $N()$  te definiëren – soms bestaan zo'n manieren misschien zelfs niet...*

*Het volgende beeld toont heel ongeveer het verschil tussen de twee definities: de ene keer ( $N_a()$ ) heb je heel veel verschillende lokale maxima en je komt elke keer in een ander lokaal maximum terecht. De andere keer ( $N_b()$ ) heb je maar één lokaal maximum dat ook het globale maximum is en waar je vanuit elk punt naartoe kan komen en waar je met elke stap vooruitgang boekt.*

*Als je het beeld van een landschap gebruikt, kan je dus zeggen dat de wereld er helemaal anders kan uitzien als je een andere definitie van  $N()$  gebruikt. Dat is natuurlijk niet echt onverwacht...*

*In de beelden die wij tekenen, is  $M$  altijd 1- of 2-dimensionaal omdat je het anders niet kan tekenen. De werkelijke dimensie is in de meeste gevallen veel groter...*



**Oefening 7** Stel dat  $G$  een complete graaf is met 30 toppen en precies één opspannende samenhangende deelgraaf  $D$  heeft minimaal gewicht. Je kiest een begingraaf door elke boog met kans  $1/2$  te kiezen. Als deze graaf niet samenhangend is, stop je. Anders pas je lokaal zoeken met de definitie  $N_a()$  van buren toe.

- Hoe groot is de kans met 50 pogingen een startgraaf te kiezen die het toelaat het optimum te vinden?
- Hoe vaak moet je het algoritme herstarten om met kans tenminste  $1/2$  tenminste één keer van een startgraaf vertrokken te zijn waar het mogelijk is  $D$  te vinden?
- Stel dat je een startgraaf met 100 bogen hebt die het toelaat  $D$  te vinden. Stel dat alle stappen gelijke kans maken en dat de graaf altijd samenhangend blijft als je maar 10 bogen verwijdert. Hoe groot is de kans dat je na 10 stappen nog altijd met een graaf zit die het toelaat  $D$  te vinden?
- Wat is de verwachtingswaarde van het aantal bogen in de startgraaf?

**Oefening 8** Toon aan, dat de definitie (b) van buren ervoor zorgt dat local search vanuit elke samenhangende startgraaf  $G$  een globaal optimum kan bereiken.

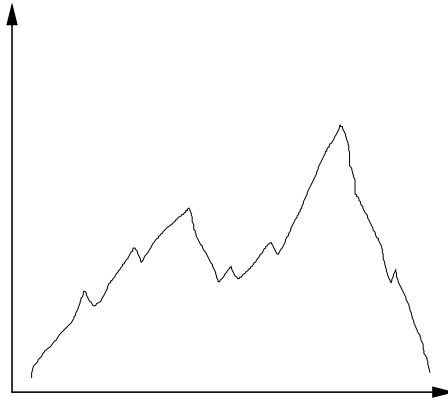
**Oefening 9** Vertaal het probleem met de minimum opspannende deelgraaf naar een ander probleem waar het lokaal zoeken algoritme het toelaat met een toevallig gekozen graaf te beginnen (ook al is hij niet samenhangend) en waar de functie die geoptimaliseerd wordt ervoor zorgt dat altijd – om het even wat de startgraaf is – de juiste grafen als optimum gevonden worden.

- Geef een definitie van de functie die geoptimaliseerd moet worden, de nieuwe definitie van  $M$  en de nieuwe definitie van buren.

## 2.2 Guided local search

Het probleem is dus dat lokaal zoeken gemakkelijk in lokale optima kan terechtkomen – zonder een kans te ontsnappen. Het onderwerp van de meeste zoekalgoritmen is vooral manieren te bedenken hoe je uit zo lokale optima kan ontsnappen – zonder echt op een toevallige manier opnieuw te beginnen. De hoop is dat dicht bij goede lokale optima er misschien zelfs nog betere zijn...

Eén mogelijkheid is het soms te aanvaarden dat je naar een slechter punt gaat. Daarvoor zullen wij een voorbeeld zien in sectie 2.4. Een andere mogelijkheid is wel altijd naar een beter punt te gaan, maar het landschap te wijzigen door een andere buurfunctie te kiezen. Dat zullen wij zien in sectie 2.3. Maar je kan het landschap ook wijzigen door een andere functie te kiezen (dus ook een andere definitie van *beter*) – en dat zullen wij nu zien.



Stel dat een verzameling  $M$  en een doelfunctie  $f()$  gegeven zijn en dat  $N()$  ook al gekozen is.

Het idee van geleid lokaal zoeken of guided local search is dat je als je in een lokaal optimum  $x$  vastzit de functie  $f()$  vervangt door een functie  $f'()$  waarvoor  $x$  geen lokaal optimum meer is en dan opnieuw lokaal zoeken toepast. Natuurlijk moet er wel een sterke samenhang tussen  $f()$  en  $f'()$  zijn. De gewone manier van doen is dat  $f()$  een deel van  $f'()$  is, maar  $f'$  zo veranderd wordt dat als je in een lokaal optimum zit een *straf* toekent aan een eigenschap van het punt  $x$ .

### **Algoritme 3 Guided local search:**

*Voorwaarde: Gegeven  $M$  en  $f()$  en wij hebben  $N()$  al gekozen. Stel dat wij het maximum zoeken.*

- a.)  $f'() = f()$  en kies een toevallig punt  $x$  en als voorlopig beste oplossing  $b$  kies je  $f(x)$ .

b.) *Herhaal de volgende stappen tot aan een zekere eindvoorwaarde is voldaan (bv. een zekere tijd is verstreken, een zeker aantal stappen is gedaan, etc.):*

- *Pas lokaal zoeken vanuit  $x$  toe om een lokaal optimum  $y$  van  $f'()$  te vinden.*
- *Als  $f(y) > b$  (**niet**  $f'(y)$ ) dan zet je  $b = f(y)$ .*
- *$x = y$*
- *Definieer een nieuwe functie  $f'()$  die het (hopelijk) toelaat het lokale optimum  $x$  te verlaten. (Hoe je dat precies kan doen, zien wij in het volgende deel.)*

c. *Geef  $b$  als resultaat terug.*

In principe klinkt dat heel gemakkelijk, maar de kunst is de functie  $f'()$  op een manier te kiezen dat de nieuwe lokale optima ook echt goed zijn voor  $f()$ . Je kan zeker vele manieren bedenken om dit idee te verwezenlijken. In het volgende gaan wij één mogelijkheid zien:

Het principe is voor sommige eigenschappen van een punt een *straf* toe te kennen. Daardoor geef je voorkeur aan oplossingen met een andere structuur – oplossingen die deze eigenschap niet hebben. Zo'n eigenschap kan van alles zijn.

- In het voorbeeld 1 bv. *je staat op een groene plek of het punt waar je staat is van steen etc.*
- In het voorbeeld 4 bv. *de maximale graad van de graaf is kleiner dan 4 of de oplossing bevat geen cykel of de oplossing bevat een cykel of er is een partitie van de toppen van  $G$  in 2 verzamelingen  $M_1, M_2$  zodat de oplossing niet de kleinste boog tussen  $M_1$  en  $M_2$  bevat of ...*

De laatste eigenschappen houden duidelijk verband met het probleem, maar de anderen kunnen inderdaad ook gebruikt worden...

Natuurlijk mogen wij alleen maar eigenschappen gebruiken, die efficiënt berekend kunnen worden – ten slotte moet een heuristiek snel zijn.

Daarbij kies je een goede mengeling van eigenschappen die lokale maxima van  $f()$  wel hebben – die kunnen je helpen om aan dergelijke maxima te ontsnappen – en die ze niet hebben – dat kan je helpen achteraf goede oplossingen voor  $f()$  terug te vinden. Wij zullen een voorbeeld uitwerken, waar wij een eigenschap **en** de tegenovergestelde eigenschap gekozen hebben en waar jullie zullen zien dat dat (in dit geval) goed werkt.

Bovendien leggen wij voor elke eigenschap  $E_i$  een gewicht  $g_i$  vast. Een groot gewicht vergroot de kans dat wij deze eigenschap vaak zullen gebruiken. Als wij nu  $k$  eigenschappen  $E_1, \dots, E_k$  hebben, dan definiëren wij de functie  $f'(x)$  als volgt:

$$f'(x) = f(x) - \lambda * \sum_{i=1}^k (a_i * I_i(x))$$

Hierbij is  $\lambda$  een getal dat wij moeten kiezen,  $a_i$  een variabele waarin wij bijhouden hoe vaak eigenschap  $i$  al in een lokaal optimum gekozen werd om  $f'()$  te wijzigen en

$$I_i(x) = \begin{cases} 1 & \text{als } x \text{ eigenschap } E_i \text{ heeft} \\ 0 & \text{anders} \end{cases}$$

Behalve  $N()$  en dit geval de  $E_i$  hebben veel metaheuristieken nog één of meerdere parameters (in dit geval zijn dat  $\lambda$ , en de  $g_i$ ) die je kan kiezen. Dan kan je testen met verschillende parameters draaien en bepalen welke parameters voor het specifieke probleem de beste zijn.

In het begin zijn dus alle  $a_i$  gelijk aan 0 en dus  $f'(x) = f(x)$ . Maar tijdens het algoritme worden de  $a_i$  gewijzigd, zodat  $f'()$  later verschilt van  $f()$ .

Als wij een optimum bereiken – dat lokaal of globaal kan zijn – berekenen wij voor elke eigenschap  $E_i$  waaraan het lokale optimum voldoet de waarde  $q_i = g_i / (1 + a_i)$  en kiezen alle  $i$  waarvoor  $q_i$  het grootst is. Deze  $a_i$ 's verhogen wij met 1 (dat is in feite het enige resultaat van het *kiezen*), waardoor wij een nieuwe functie  $f'()$  krijgen die het misschien toelaat dit optimum te verlaten en een nieuwe te vinden.

**Voorbeeld 5** *Wij gebruiken het voorbeeld 4 met de definitie  $N_a$  van burens. Wij weten al dat je zo heel gemakkelijk in een lokaal optimum terecht kan komen. Wij gebruiken de eigenschappen*

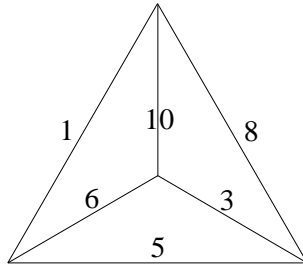
$E_1$  *De maximale graad van de graaf is kleiner dan 4.*

$E_2$  *De oplossing bevat geen cykel.*

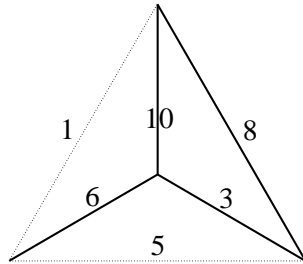
$E_3$  *De oplossing bevat een cykel.*

*Wij kiezen  $\lambda = 2$  en  $g_1 = 2, g_2 = g_3 = 1$ .  $f(x)$  is de som van de gewichten van de bogen in de graaf  $x$  en wij zoeken het minimum. (Wij zouden natuurlijk even goed het maximum kunnen zoeken nadat we de functie met  $-1$  vermenigvuldigd hebben.)*

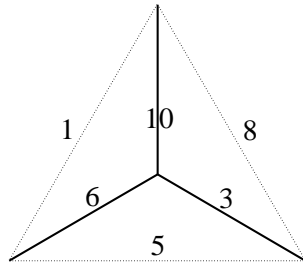
*Stel dat de graaf  $G$  de  $K_4$  is met gewichten als volgt:*



en dat de random startgraaf  $x_0$  zoals volgt is:



Dus  $f(x_0) = 27$ . Als wij nu lokaal zoeken toepassen, dan zou bv. het verwijderen van de boog met gewicht 8 een bewerking zijn, die aanvaard wordt en ons naar een lokaal optimum leidt waar lokaal zoeken stopt:



Voor dit optimum is  $f(x) = 19$  – het is dus een nieuwe beste oplossing. Deze graaf heeft eigenschappen  $E_1$  en  $E_2$  en wij berekenen  $q_1 = 2/(1+0) = 2$  en  $q_2 = 1/(1+0) = 1$ . Dus wordt  $E_1$  gekozen en wij hebben nu  $a_1 = 1$  en (let op – wij zoeken hier een minimum in plaats van een maximum)

$$f'(x) = \begin{cases} f(x) + 2 & \text{als alle graden kleiner dan 4 zijn} \\ f(x) & \text{anders} \end{cases}$$

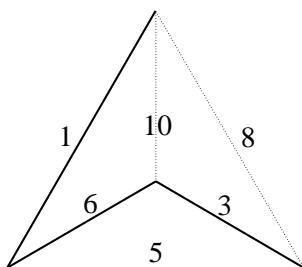
Als wij opnieuw lokaal zoeken toepassen, zien wij, dat wij het lokale minimum niet kunnen verlaten – het nieuwe lokale minimum is het oude.

Maar nu is  $q_1 = 2/(1+1) = 1$  en  $q_2 = 1/(1+0) = 1$ . Dus worden  $E_1$  en  $E_2$  gekozen en wij hebben nu  $a_1 = 2$  en  $a_2 = 1$  en

$$f'(x) = f(x) + e_1 + e_2$$

waarbij  $e_1 = 4$  als  $x$  eigenschap  $E_1$  heeft en anders 0 en  $e_2 = 2$  als  $x$  eigenschap  $E_2$  heeft en anders 0.

Het lokale optimum  $x$  heeft dus  $f'(x) = 19 + 4 + 2 = 25$ . Als wij opnieuw lokaal zoeken toepassen zien wij dat wij nu wel kunnen ontsnappen: Door de boog met gewicht 1 toe te voegen krijgen wij een cykel in de nieuwe graaf  $x'$  en dus  $f'(x') = 20 + 4 = 24$ . Daarna kan lokaal zoeken bv. boog 10 verwijderen. dan hebben wij  $f'(x'') = 10 + 4 + 2 = 16$  en wij zitten in een nieuw lokaal minimum waarvoor de functie met  $f(x) = 10$  een nieuwe beste waarde oplevert.



Dat is nog altijd niet het globale optimum maar het is duidelijk hoe het werkt. En wij hebben ook gezien waarom het goed was niet gewoon van een nieuwe startgraaf te vertrekken: een beetje informatie over het globale optimum zit ook in de lokale optima – in dit geval bevatten ook lokale optima vaak veel bogen met een klein gewicht.

**Oefening 10** • Kan je als je met het algoritme doorgaat het globale optimum ook vinden?

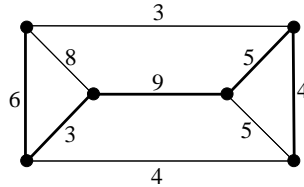
- Zijn de waarden van  $\lambda$  en de  $g_i$  goed gekozen? Welke waarden zouden beter zijn?

Natuurlijk is dit een klein en onnozel voorbeeld maar jammer genoeg zijn echte en interessante voorbeelden niet met de hand te doen ...

Inderdaad kunnen de ideeën van *guided local search* natuurlijk niet alleen op lokaal zoeken toegepast worden – ze werken met om het even welke metaheuristiek je gebruikt om optima te vinden. Je moet gewoon de stap waar lokaal zoeken wordt gebruikt door een stap vervangen waar je een andere metaheuristiek gebruikt. Maar omdat lokaal zoeken gemakkelijk in slechte lokale optima terechtkomt terwijl de andere metaheuristicen zelf ideeën bevatten om dat te voorkomen is het voor lokaal zoeken bijzonder belangrijk.



**Oefening 11** Pas het beschreven algoritme voor *guided local search* toe op de volgende graaf met de gegeven startdeelgraaf:



Als je betere eigenschappen kan bedenken of betere parameters kan kiezen dan die in de les, gebruik die dan!

**Oefening 12** De volgende oefening was een deel van een examen:

- Beschrijf expliciet een *guided local search* algoritme voor het volgende probleem. Beschrijf alle definities en parameters die nodig zijn om jouw algoritme expliciet te beschrijven. Maar jouw parameters en de details van het algoritme moeten natuurlijk niet optimaal gekozen zijn – de bedoeling is dat je toont dat je *guided local search* verstaan hebt en op problemen kan toepassen.

**Probleem:**

Een supermarkt heeft verschillende artikelen  $a_1, \dots, a_n$  die onmiddellijk naar een andere supermarkt vervoerd moeten worden. De artikelen behoren tot verschillende klassen, zoals bv. *electronica*, *sportartikelen*, *kledij*, etc. Elk artikel  $a_i$  heeft een waarde  $w(a_i)$  en een gewicht  $g(a_i)$ . Jammer genoeg staat er maar één vrachtwagen ter beschikking en die kan maar een maximaal gewicht van  $g$  dragen. Het doel is nu een verzameling van artikelen samen te stellen zodat de som van de gewichten ten hoogste  $g$  en de som van de waarden maximaal is.

## 2.3 Variable neighbourhood metaheuristieken

Terwijl *guided local search* de definitie van de doelfunctie wijzigt om uit een lokaal optimum te ontsnappen, wijzigt – of beter: vervangt – *variable neighbourhood search* de definitie van buur (dus  $N()$ ) om te ontsnappen. Een globaal optimum is natuurlijk een optimum voor **elke** definitie van  $N()$ . Als elk punt weinig burens heeft, kan je sneller zoeken en als er een buur is die naar een goede oplossing leidt, is de kans bovendien groter dat je hem ook kiest (zie Oefening 6). Daarom is het goed als een punt niet te veel burens heeft. Jammer genoeg kan je dan snel in een lokaal optimum vastzitten.

Variable neighbourhood descent en variable neighbourhood search gebruiken een compromis: zij gebruiken meerdere definities van burens – maar niet tegelijk. Op deze manier heb je meestal relatief weinig burens die je moet onderzoeken en aan de andere kant als je naar alle definities kijkt toch veel richtingen die je in kan gaan om uit slechts lokale optima te ontsnappen.. Beide metaheuristieken gebruiken de ene definitie met niet te veel burens per punt om snel een lokaal optimum te vinden en dan een andere definitie om (hopelijk) uit dit lokale optimum te ontsnappen.

#### Algoritme 4 Variable neighbourhood descent

*Gegeven een verzameling  $M$ , een doelfunctie  $f : M \rightarrow \mathbb{R}$  waarvan wij het maximum zoeken. Wij hebben  $m$  buurfuncties  $N_1(), \dots, N_m()$  gekozen. Bovendien hebben wij misschien een voorwaarde om te stoppen (tijd, aantal stappen, etc.).*

(a) *Kies een startpunt  $x$  – bv. op een toevallige manier.*

(b) *Herhaal de volgende stappen tot aan een stopvoorwaarde is voldaan”:*

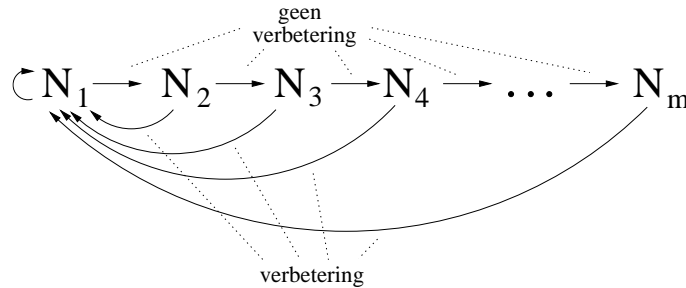
- *begin met  $k = 1$  en definitie  $N_k() = N_1()$*
- *herhaal de volgende stappen tot  $k = m + 1$ :*
  - Zoeken:** *Zoek het beste punt  $x' \in N_k(x)$  (dus  $x'$  met  $f(x') \geq f(x'') \forall x'' \in N_k(x)$ ).*
  - Volgende stap vastleggen:**
    - *Als  $f(x') > f(x)$  neem dan  $x = x'$  en begin opnieuw met  $k = 1$ .*
    - *Anders:  $k = k + 1$ . Als  $k = m + 1$ : STOP.*

Het is dus gewoon een gretige versie van lokaal zoeken met verschillende definities van buur en je stopt zodra geen enkele definitie van buur naar een beter punt leidt. Het feit dat je naar alle burens moet kijken zorgt er al voor dat elke functie  $N()$  niet te veel burens per punt mag vastleggen om efficiënt te zijn.

Het *descent* in VND is natuurlijk alleen maar juist als je een minimum zoekt. Het is duidelijk wat je daarvoor moet wijzigen en dan is het inderdaad een afdaling.

**Oefening 13** *Kan de prestatie van VND veranderen als de volgorde van de functies  $N_1(), \dots, N_m()$  verandert?*

Als je ernaar kijkt welke buurfunctie in de lus van Algoritme 4 gebruikt wordt dan kan je dat ongeveer door het volgende stroomdiagram voorstellen:



Terwijl dus  $N_1$  voor elk punt gebruikt wordt (dus heel vaak), wordt  $N_m$  alleen dan gebruikt als in geen van de andere omgevingen een beter punt gevonden kan worden. Als de omgevingen dus goed gekozen zijn, zal  $N_m$  niet vaak doorzocht moeten worden. Omdat in de lus van Algoritme 4 de hele omgeving doorzocht moet worden, is het dus bijzonder belangrijk voor de efficiëntie dat  $N_1$  klein is – terwijl  $N_m$  duidelijk groter mag zijn omdat deze omgeving niet vaak doorzocht wordt.

**Oefening 14** *Metaheuristieken zijn vooral belangrijk voor heel moeilijke problemen – bv. NP-complete problemen. Stel dus dat  $P \neq NP$  wat betekent dat voor deze algoritmen geen polynomiale algoritmen bestaan. Wie niet meer weet wat NP-compleet betekent, zou dat voor deze oefening beter nog eens opzoeken...*

Gegeven een graaf  $G_0$ . Zij dan  $M = \{G | G \text{ is deelgraaf van } G_0\}$ . Wij schrijven  $h(G)$  voor de lengte van een kortste cykel in  $G$ . Als er geen cyclen zijn, is  $h(G) = 0$ . De doelfunctie is  $f(G) = 2 * h(G) - |E(G)|$  en wij zoeken het maximum.

- Kan je een deelgraaf beschrijven waarvoor deze functie maximaal is?

Wij willen variable neighbourhood descent toepassen voor een vast aantal  $k$  van buurfuncties  $N_1(), \dots, N_k()$ . Stel dat de buurfuncties de burens in een zekere volgorde geven en niet als een verzameling zonder orde. Als er in de zoek-stap van VND meerdere even goede punten (grafen) zijn waar je naartoe kan gaan, neem je altijd de eerste van deze grafen in de volgorde van de buurfunctie die je gebruikt.

Het is belangrijk voor de efficiëntie dat je niet te veel burens hebt en dat die ook niet te ingewikkeld berekend moeten worden. Stel dus dat voor elke graaf en elke  $N_i()$  geldt  $|N_i(G)| \leq |E(G_0)|^4$  en dat het ten hoogste lineaire tijd per buur vraagt om de burens te berekenen.

Natuurlijk proberen wij de  $N_i()$ 's op een manier te kiezen dat VND zo goed mogelijk presteert – maar hoe goed kan dat ?

- Toon aan dat het (met de gegeven voorwaarden) niet mogelijk is  $N_1(), \dots, N_k()$  te vinden zodat het mogelijk is dat VND vanuit elke begingraaf het optimum kan bereiken.
- Is het mogelijk  $N_1(), \dots, N_k()$  te vinden zo dat VND als het vanuit  $G_0$  vertrekt het optimum kan vinden?
- Werken jouw argumenten ook als je bv. in het geval van meerdere even goede punten voor de volgende stap de burens toevallig kiest?
- Kan je iets erover zeggen hoe moeilijk het is voor een gegeven  $G_0$  goede begingrafen te berekenen – dat zijn begingrafen zodat je met VND altijd het optimum vindt?

**Oefening 15** Geef 2 functies  $N_1(), N_2()$  voor het voorbeeld 4 zo dat voor geen van de twee functies lokaal zoeken altijd het optimum vindt – maar VND wel.

Deze VND-heuristiek kan je alleenstaand gebruiken of als een deel van *variable neighbourhood search*:

#### Algoritme 5 Variable neighbourhood search

Gegeven een verzameling  $M$ , een doelfunctie  $f()$  waarvan wij het maximum zoeken. Wij hebben  $m$  buurfuncties  $N_1(), \dots, N_m()$  gekozen. Bovendien hebben wij een voorwaarde om te stoppen (tijd, aantal stappen, etc.).

(a) Kies een startpunt  $x$  – bv. op een toevallige manier.

(b) Herhaal de volgende stappen tot aan de stopvoorwaarde is voldaan:

- begin met  $k = 1$  en buurfunctie  $N_k = N_1$
- herhaal de volgende stappen tot en met  $k = m$

**Schudden:** Kies een toevallig punt  $x' \in N_k(x)$ . Houd er geen rekening mee of  $f(x') \geq f(x)$ .

**Zoeken:** Zoek een lokaal optimum  $x''$  door van het net gekozen punt  $x'$  te vertrekken. Dat kan je doen door lokaal zoeken waarbij definitie  $N_1$  of  $N_k$  wordt toegepast, of beter door *variable neighbourhood descent* toe te passen.

**Vergelijken:** – Als het nieuwe optimum  $x''$  beter is dan het tot nu toe beste  $x$ , kies  $x = x''$  en  $k = 1$ .

– Anders verhoog  $k$  dus  $k \leftarrow k + 1$ .

**Oefening 16** *Ontwerp een VNS algoritme voor het handelsreizigersprobleem. Gebruik ten minste 3 functies  $N_1()$ ,  $N_2()$ ,  $N_3()$ . Pas jouw algoritme toe op een klein voorbeeldje.*

**Oefening 17** *Gegeven een graaf.  $M$  zij de verzameling van alle deelgrafen van  $G$  en voor  $G \in M$  zij  $f(G) = 1$  als  $G$  een opspannende cykel (dus een Hamiltoniaanse cykel) is en  $f(G) = 0$  anders. Van welke van de tot nu toe gezien metaheuristieken denk je dat die het best presteert. Wat is het probleem met deze functie? Hoe kan je dat oplossen?*

Volgens Pierre Hansen, één van de uitvinders van VNS, presteren buurfuncties met de eigenschap dat voor elk punt  $x$  geldt  $N_1(x) \subset N_2(x) \subset \dots \subset N_m(x)$  vaak heel goed. Dat is het geval als in de omgeving (volgens de definitie  $N_1()$ ) van een lokaal optimum nog meer en betere optima zitten – misschien alleen maar iets verder weg dan de directe omgeving. Natuurlijk moet je als je VND met deze buurfuncties toepast en het beste punt in  $N_i()$  zoekt het deel dat in  $N_{i-1}()$  ligt niet opnieuw testen.

Maar je kan de andere  $N_i()$  ook op een manier kiezen dat het landschap helemaal verandert – dus waar bv.  $N_2()$  niet zo direct iets te maken heeft met  $N_1()$ .

Als wij het idee (of de hoop) gebruiken dat er in de buurt van een optimum vaak nog betere optima zijn, dan is het dus zinvol  $N_1()$  relatief klein te kiezen om ook een grotere kans te hebben een punt te vinden dat als startpunt op weg naar dit betere lokale optimum kan dienen (vergelijk Oefening 6). Pas later als met deze kleine omgevingen het doel niet wordt bereikt, worden ook grotere omgevingen gebruikt.

Een heel informele beschrijving van het voordeel van VNS in vergelijking met VND is dat terwijl VND echt een beter punt in een omgeving moet vinden om naar het volgende optimum te kunnen gaan, het voor VNS al voldoende is het begin van de helling te vinden (maar dat is natuurlijk **heel** informeel. . .)

## 2.4 Simulated annealing

Behalve in één stap in VNS zijn wij tot nu toe altijd naar betere punten gegaan. Wij hebben andere doelfuncties gekozen om uit lokale optima te kunnen ontsnappen, wij hebben de definities van buur gewijzigd, maar wij zijn nooit naar punten gegaan waar de op dat moment gebruikte doelfunctie slechter was.

Als je opnieuw het beeld van een landschap en het hoogste punt gebruikt dan lijkt het vrij kunstmatig het landschap (de definitie van buur) of de interpretatie van *hoogte* (de doelfunctie) te wijzigen. De enige manier om

uit een lokaal minimum te ontsnappen is gewoon naar beneden te gaan – of met andere woorden: soms ook naar slechtere punten te gaan. De vraag is natuurlijk **wanneer** je naar slechtere punten gaat. Als je altijd naar punten in de buurt gaat – om het even of ze beter of slechter zijn – loop je gewoon toevallig rond. Dan kan je beter helemaal toevallige punten kiezen (dus *toevallig zoeken*).

In S. Kirkpatrick, C.D. Gelatt, M.P. Vecchi: *Optimization by Simulated Annealing*, Science 220, n. 4598, pp 671–680 (1983) werd *simulated annealing* voor de eerste keer voorgesteld.

De techniek van simulated annealing of gesimuleerd temperen is inderdaad het nabootsen van wat mensen deden al lang voordat computers werden uitgevonden. Als ijzer vloeibaar is, hebben de atomen een heel grote energie – voldoende om niet op een manier gebonden te zijn die optimaal is qua energie. Ze kunnen dan ook naar slechtere posities gaan – gewoon omdat ze ervoor voldoende energie hebben. Tijdens het afkoelen wordt het dan altijd moeilijker energetisch slechte posities te bereiken en ten slotte zit elk atoom in een lokaal minimum waaruit het niet kan ontsnappen. Hoe dicht dat bij een globaal minimum zit, hangt van meerdere parameters af – bv. hoe snel het materiaal afkoelt. Voor sommige materialen is het temperen ook nadat ze gelast zijn heel belangrijk om een stevig en duurzaam materiaal te hebben. De exacte parameters (hoe heet in het begin, hoe snel afkoelen, etc.) zijn voor elk materiaal verschillend. En inderdaad zijn de corresponderende parameters ook voor **gesimuleerd** temperen heel belangrijk.

Het idee is dus dat je in het begin een vrij hoge temperatuur hebt. Dat stelt je in staat ook naar punten te gaan die (qua energie – of in ons geval qua doelfunctie) veel slechter zijn dan waar je bent. Of precies: de **kans** is groot dat het lukt. Maar hoe meer tijd verstrijkt (hoe meer stappen je doet), hoe minder hoog de temperatuur is – dus hoe minder je naar slechtere punten kan gaan (of precies: hoe minder de **kans** is dat dat lukt). Inderdaad is het in de realiteit zo (of tenminste in het model dat de fysica ervan heeft) dat het voor een zekere temperatuur niet zo is, dat sommige stappen naar een slechter punt altijd aanvaard worden en andere nooit, maar het is zo, dat alle met een zekere **kans** aanvaard worden. De kans van  $x$  naar  $x'$  te gaan daalt gewoon als  $f(x) - f(x')$  groot is (als wij het maximum zoeken). En bovendien daalt de kans ook als de temperatuur kleiner is. Dus is het *in principe* nog altijd mogelijk ook naar een duidelijk slechter punt te gaan als de temperatuur klein is – de kans dat echt te doen is gewoon heel klein.

Er zijn ook heuristieken, die geen kansen gebruiken, maar met drempels werken (thresholds) – dus tot een zekere drempel slechtere punten **altijd** aanvaarden en boven de drempel **nooit**. Ook dergelijke metaheuristieken (*threshold accepting* of *great deluge*) werken heel goed in sommige gevallen.

Soms beter dan simulated annealing, soms slechter... Maar dat is iets voor een ander hoofdstuk (niet in deze les).

Dus wat hebben wij nodig?

- Een temperatuurfunctie  $t : \mathbb{N} \rightarrow \mathbb{R}_+ \setminus \{0\}$ . De waarde  $t(k)$  zegt hoe hoog de temperatuur in stap  $k$  is. Om goed te werken zou  $t()$  het best nooit groeien en  $\lim_{k \rightarrow \infty} t(k) = 0$ .
- Een functie  $p : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$  die bepaalt hoe groot de kans is, een gekozen stap ook echt te doen. Als je in punt  $x$  bent en naar punt  $x'$  wilt gaan, waarbij  $\Delta = f(x') - f(x)$  en de temperatuur is  $t$ , dan is  $p(\Delta, t)$  de kans dat je dat ook echt doet. De kans hangt dus niet direct van de punten af, maar alleen van het verschil in de doelfuncties, maar je kan de functie natuurlijk ook schrijven als  $p' : M \times M \times \mathbb{R} \rightarrow [0, 1]$  met  $p'(x', x, t) = p(f(x') - f(x), t)$ .

Een mogelijke functie  $t()$  is bijvoorbeeld  $t(n) = T_0 * \lambda^{\lfloor n/L \rfloor}$  waarbij de begintemperatuur  $T_0 \in \mathbb{R}$ ,  $T_0 > 0$  en  $L \in \mathbb{N}$ ,  $L > 0$  en  $\lambda \in \mathbb{R}$ ,  $0 < \lambda < 1$  parameters zijn, die voor elke toepassing gekozen moeten worden – bv. door testen te draaien met verschillende waarden en de resultaten te vergelijken. Eén manier om  $T_0$  te kiezen is bv. een bovengrens te bepalen voor het maximale verschil  $|f(x') - f(x)|$  als dat gemakkelijk is om te bepalen. Dat werkt normaal vrij goed maar je kan  $T_0$  ook anders kiezen.

Als het maximum wordt gezocht – het dus in principe goed is naar  $x'$  te gaan als  $f(x') > f(x)$  wordt normaal als  $p()$  de volgende functie gebruikt:

$$p(\Delta, t) = \begin{cases} 1 & \text{als } \Delta > 0 \text{ (dus } f(x') > f(x)) \\ e^{\Delta/t} & \text{anders. Omdat } \Delta \text{ negatief of nul is, is dat } \leq 1 \end{cases}$$

Deze functie beschrijft ongeveer de realiteit van het temperen van metalen – ze komt dus uit de fysica en kan met hulp van de Boltzmann verdeling berekend worden. Algoritmen die deze functie ook voor simulated annealing toepassen, werken vaak vrij goed – maar *in principe* kan je natuurlijk ook andere functies gebruiken die ongeveer vergelijkbare eigenschappen hebben. Het feit dat ook algoritmen die de metaheuristiek *threshold accepting* toepassen goed presteren, toont dat de details van deze functie niet altijd belangrijk zijn. Je kan threshold accepting interpreteren als simulated annealing met een heel gemakkelijke functie  $p()$ .

**Oefening 18** Wat moet je aan  $p()$  wijzigen als je het minimum zoekt in plaats van het maximum?

**Oefening 19** Welke eigenschappen zou jij van zo'n functie  $p()$  eisen om te kunnen verwachten dat jouw algoritme goed werkt?

Met welke eigenschappen zou jij zeker zijn dat het algoritme niet goed werkt?

### Algoritme 6 Simulated annealing

Gegeven een verzameling  $M$ , een doelfunctie  $f()$  waarvan wij het maximum zoeken, en wij hebben de buurfunctie  $N()$  en de functies  $t()$  en  $p()$  al gekozen. Bovendien hebben wij een voorwaarde om te stoppen (tijd, aantal stappen, etc.).

(a) Kies een startpunt  $x$  – bv. op een toevallige manier.

(b) Herhaal de volgende stappen tot aan de stopvoorwaarde is voldaan. Begin met  $s = 0$  ( $s$  houdt het aantal stappen bij)

- Kies een toevallig punt  $x' \in N(x)$ .
- Bereken  $\Delta = f(x') - f(x)$ .
- Kies  $x = x'$  met kans  $p(\Delta, t(s))$ .
- $s = s + 1$

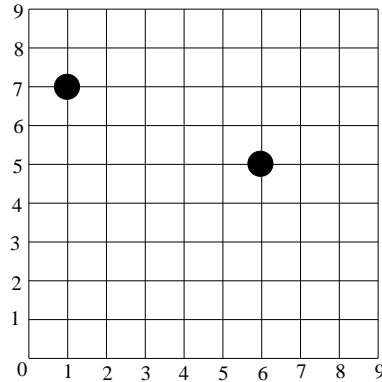
(c) Geef  $x$  als oplossing.

Omdat simulated annealing een relatief oude techniek is, bestaan er natuurlijk al heel veel artikels over hoe je het basisalgoritme kan verbeteren, hoe je de parameters voor sommige problemen moet kiezen, etc.

Iets dat onmiddellijk een goed idee lijkt is bv. het beste gevonden punt altijd bij te houden. Inderdaad werd aangetoond dat simulated annealing (onder zekere voorwaarden) naar het optimum convergeert – maar als je dat wilt garanderen dan vraagt het algoritme in de praktijk meer tijd dan je ervoor kan/wil besteden.

**Oefening 20** De bedoeling is 2 verschillende toppen  $p_1, p_2$  van een  $10 \times 10$  tralie met een zekere eigenschap te vinden: Als de toppen de coördinaten  $(x, y)$  en  $(x', y')$  hebben dan is de afstand  $d((x, y), (x', y')) = |x - x'| + |y - y'|$ . Wij hebben  $M = \{(x, y, x', y') | 0 \leq x, y, x', y' \leq 9, (x, y) \neq (x', y')\}$  dus de verzameling van alle mogelijke posities van de twee punten. De doelfunctie is gedefinieerd als  $f(x, y, x', y') = |2 - d((4, 5), (x, y))| + |3 - d((x, y), (x', y'))| + |2 - d((5, 5), (x', y'))|$ . Wij zoeken het minimum van  $f()$ .





Definieer functies  $N()$ ,  $t()$  en  $p()$  en pas simulated annealing toe door sommige stappen met de hand uit te werken.

**Oefening 21** Gegeven  $k$  punten  $p_1, \dots, p_k$  in het vlak (dus in  $\mathbb{R}^2$ ). Ze moeten op een manier geplaatst worden dat voor elk paar geldt dat de afstand zo dicht mogelijk bij 1 is. Natuurlijk is het voor grote  $k$  niet mogelijk dat op een perfecte manier te doen.

Wij definiëren de fout als

$$f(p_1, \dots, p_k) = \sum_{1 \leq i < j \leq k} (d(p_i, p_j) - 1)^2$$

waarbij  $d()$  de gewone Euclidische afstand is.

- Wat is de grootste  $k$  waarvoor je een oplossing met fout 0 kan vinden? Toon aan dat jouw antwoord juist is.
- Ontwerp een simulated annealing algoritme dat vertrekt van een toevallige plaatsing van de punten in het vierkant  $[0, 10] \times [0, 10]$ . Denk er goed over na wat je als mogelijke stappen wil kiezen en voorspel wat jouw algoritme gaat doen. Houd rekening met het aantal mogelijkheden voor elke stap (degrees of freedom).
- Implementeer jouw algoritme en vergelijk het met dat van jouw collega's. Kloppen jouw voorspellingen? Zijn er programma's die duidelijk beter/slechter presteren? Waarom?

**Oefening 22** Gegeven  $n$  reële getallen  $0 < a_1, \dots, a_n < 1$ . Gezocht is een deelverzameling  $a_{i_1}, \dots, a_{i_k}$  van getallen zodat

$$0 < \sum_{j=1}^k a_{i_j} \leq 1$$

en deze som maximaal is.

Ontwerp een *simulated annealing* algoritme voor dit probleem.

Geef alle definities en parameters die nodig zijn om het algoritme volledig te beschrijven maar het is niet nodig dat ze ook zo gekozen zijn dat het algoritme goed werkt. Daarvoor zou je het echt moeten implementeren en testen laten draaien. Het algoritme moet alleen maar in principe werken en het optimum kunnen vinden (ook al zou dat veel te lang duren en te veel pogingen vragen). De bedoeling is gewoon om aan te tonen dat het principe verstaan is.

## 2.5 Genetische algoritmen

Zoals *simulated annealing* proberen ook genetische algoritmen een optimalisatieproces dat al lang gekend is op de computer na te bootsen. In dit geval is dat de evolutie. De doelfunctie die door de evolutie geoptimaliseerd wordt is de fitness – dat is hoe goed een individu bij zijn omgeving past. Inderdaad is de *fitness* niet altijd **precies** de doelfunctie – soms is het noodzakelijk de doelfunctie een beetje te wijzigen (dat gaan wij ook in een voorbeeld zien), maar *in principe* bepaalt de doelfunctie de fitnessfunctie. Iets dat je genetische algoritmen zou kunnen noemen werd al in de jaren 60 voorgesteld (Rechenberg, Schwefel) toen de computers er nog helemaal niet klaar voor waren. Maar toen was het vooral het principe van *mutatie en selectie* van één individu dat werd gesimuleerd. Het resultaat is iets dat je goed met lokaal zoeken kan vergelijken: mutatie betekent een stap naar een buur en selectie betekent deze stap aanvaarden of niet aanvaarden.

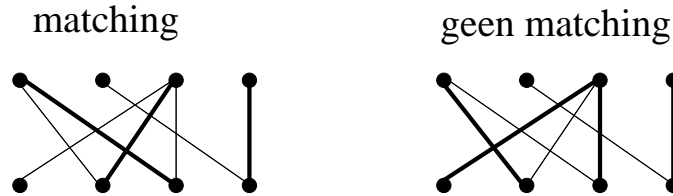
Op dit moment wordt het begrip *genetische algoritmen* vooral gebruikt voor algoritmen van het soort die wij gaan zien. Deze algoritmen werden in de jaren 70 geïntroduceerd (Holland 1975) en simuleren de evolutie van een hele populatie.

Als je de evolutie **heel** kort wil beschrijven, dan kan je dat als volgt: je begint met een populatie en elk individu heeft een zekere fitness. In elke stap worden kinderen verwekt en sommige individuen sterven. De kans kinderen te verwekken en de kans te sneuvelen hangen allebei af van de fitness van de individuen. Als er kinderen verwekt worden dan zijn dat geen identieke kopieën van hun ouders, maar er gebeuren kleine veranderingen (mutaties) en de eigenschappen van de ouders worden ook gemengd (crossover). Na voldoende stappen is de fitness van de hele bevolking gestegen. . . Dit proces gaan wij simuleren.

Voor genetische algoritmen gaan wij een nieuw voorbeeld gebruiken:

**Voorbeeld 6** Gegeven een graaf  $G = (V, E)$ . Een *matching* is een verzameling  $E' \subseteq E$  zodat elke top in ten hoogste één boog uit  $E'$  zit. Ons doel is

voor een gegeven graaf  $G$  de grootste matching te bepalen.



Ook voor dit probleem zou je het best geen metaheuristiek toepassen – er bestaan polynomiale algoritmen – maar het is heel geschikt om de principes van GA te tonen.

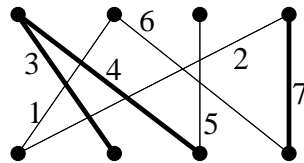
Als verzameling  $M$  gaan wij alle mogelijke verzamelingen van bogen van  $G$  kiezen – dus  $M = \{E' | E' \subseteq E\}$ . Dan kunnen wij als fitness natuurlijk niet het aantal bogen kiezen omdat wij dan altijd lokale extrema zouden vinden die geen matchings zijn. Voor een verzameling  $E'$  van bogen definiëren wij  $f(E') = |E'| - 2 * |E'_f|$  waarbij  $E'_f$  de deelverzameling van bogen van  $E'$  is waarvan tenminste één eindpunt in tenminste 2 bogen zit (dus  $E'_f = \{e \in E' | \exists e' \in E', e' \neq e, e \cap e' \neq \emptyset\}$ ). Het is gemakkelijk om in te zien dat het globale maximum altijd een deelverzameling van bogen is die een maximale matching vormt. Als wij één van de metaheuristieken zouden toepassen waar een buurfunctie  $N()$  gedefinieerd is – waar wij dus ook lokale maxima hebben – zouden bij de meeste keuzes van buurfuncties ook alle lokale maxima matchings zijn.

Voor genetische algoritmen (wij schrijven kort GA) wordt in de meeste gevallen een expliciete voorstelling van een punt in  $M$  als een string met tekens uit een eindig alfabet (b.v  $\{0, 1\}$  of zoals in de biologie  $\{A, C, G, T\}$ ) geëist. Deze string heet het genotype en het punt dat het voorstelt het fenotype – zoals ook in de biologie. Het fenotype is dus waar wij echt in geïnteresseerd zijn en het genotype zo iets als het gen van een mens – een codering.

Wij hebben dus een functie  $c : M \rightarrow A^*$  van de verzameling  $M$  naar een verzameling  $A^*$  van eindige strings uit een alfabet  $A$ . Daarbij is  $c(x)$  de string die het punt  $x$  voorstelt. Later gaan wij met de strings werken en dan hebben wij een ander probleem: Als wij de fitness van een string willen berekenen dan moeten wij weten welk punt een gegeven string voorstelt om de fitness van dit punt te berekenen. De fitness van een string  $s$  (die tot een punt in  $M$  behoort) is gedefinieerd als  $f(c^{-1}(s))$ . Wij moeten dus  $c^{-1}()$  berekenen – maar het kan gebeuren dat sommige strings  $s$  helemaal geen punten uit  $M$  voorstellen, dus dat  $c^{-1}(s)$  niet gedefinieerd is. Als  $M$  bijvoorbeeld de verzameling van alle mensen is en  $S$  de verzameling van alle eindige strings uit de letters A,C,G,T (dus  $S = \{A, C, G, T\}^*$ ). Dan heeft elke mens een voorstelling in  $S$ , maar de

meeste strings in  $S$  behoren tot geen enkele mens (of dier of wat dan ook). Eén van de grootste problemen met het ontwerp van een GA is de codering en de veranderingen (mutatie, crossover) op een manier te kiezen dat als je iets aan een string verandert het resultaat (tenminste meestal) ook tot een punt behoort. Maar dat is veel gemakkelijker gezegd dan gedaan...

In ons voorbeeld coderen wij de punten (verzamelingen van bogen) als volgt: wij geven nummers  $1, \dots, |E|$  aan de bogen in  $G$ . Als  $e_i$  de boog met nummer  $i$  is, dan is dus  $E = \{e_1, \dots, e_{|E|}\}$ . Wij kiezen als alfabet  $\{0, 1\}$  en stellen de verzamelingen voor als 0,1-vectoren van lengte  $|E|$ . Als  $E' = \{e_{i_1}, \dots, e_{i_m}\}$  dan is  $c(E')$  de 0,1-vector waarbij de  $k$ -de component 1 is als  $k \in \{i_1, \dots, i_m\}$  en 0 anders.



$E'$  met dikke lijnen

$$c(E') = (0, 0, 1, 1, 0, 0, 1)$$

$$f(E') = 3 - 2 * 2 = -1$$

Stel dat wij een codering van de punten als een string hebben – dus een genotype. Terug naar het algemeen geval.

**Mutatie:** Een mutatie is in principe hetzelfde als een stap van een punt naar een buur. Wij beschrijven mutaties op het niveau van het genotype, dus de string: een mutatie kan in veel gevallen geïmplementeerd worden als het vervangen van één of meerdere elementen in de string door een ander element uit het alfabet. Normaal wordt op een gerandomiseerde manier erover beslist welke plaatsen vervangen worden en waardoor. Het probleem en de codering  $c()$  bepalen of alle vervangingen toegelaten moeten worden of alleen maar sommige vervangingen zinvol zijn.

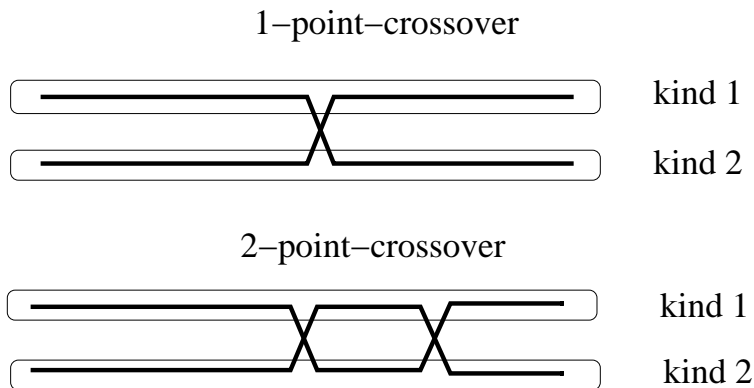
Je kan als *mutatie* ook tekens van volgorde verwisselen of andere wijzigingen op de genotype toepassen. Belangrijk is dat het een toevallige wijziging van één individu is.

Hoe groot de kans gekozen moet worden om een element van de string door mutatie te veranderen, is voor elk probleem verschillend en kan het best bepaald worden door testen te draaien. Als de kans te groot is dan zou het zo zijn alsof je op een toevallige manier een ander individu kiest – je verliest dan alle eigenschappen van het individu en dat mag zeker niet.

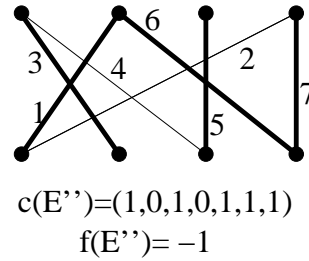
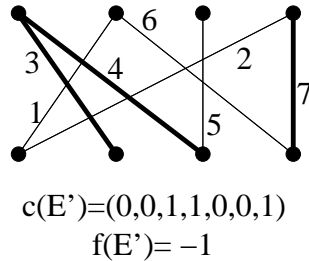
In ons voorbeeld zouden wij bv. voor een string  $(s_1, \dots, s_7)$  voor elke  $i \in \{1, \dots, 7\}$   $s_i$  met een kans van  $1/10$  vervangen door  $1 - s_i$ . Als wij de wijzigingen met kans  $1/2$  zouden doen, zou het gewoon een toevallige matching zijn – een zeer slechte keuze.

Hier zien wij al dat het moeilijk geweest zou zijn als wij als  $M$  alleen de verzameling van matchings gekozen zouden hebben: als wij in onze codering een 0 door een 1 vervangen, is de kans groot dat het resultaat geen matching meer is – de code dus geen punt codeert.

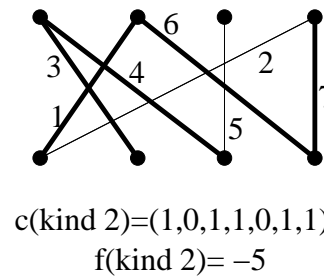
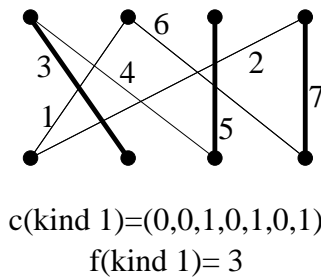
**Crossover:** Crossover gebeurt als 2 individuen één of meerdere kinderen verwekken. Het principe is dat het kind of de kinderen combinaties zijn van kopieën van de ouders. Er zijn meerdere manieren om dat te doen. De standaardmanieren zijn 1-point-crossover, 2-point-crossover, etc. Als wij b.v. 1-point-crossover willen toepassen, kunnen wij (bv.) op een random manier een index  $i$  kiezen en het eerste kind krijgt de eerste  $i - 1$  plaatsen van ouder 1 en de rest van ouder 2. Het tweede kind krijgt de eerste  $i - 1$  plaatsen van ouder 2 en de rest van ouder 1. Dus: als de ouders  $s_1, \dots, s_k$  en  $t_1, \dots, t_k$  zijn, dan zijn de kinderen  $s_1, \dots, s_{i-1}, t_i, \dots, t_k$  en  $t_1, \dots, t_{i-1}, s_i, \dots, s_k$ . Zoals veel van de details hangt de keuze of 1-point-crossover of multipoint-crossover beter is af van het probleem. Ook de manier waarop je de crossover punten kiest heeft een sterke invloed op de performantie.



In ons voorbeeld zou een 2-point-crossover met 2 kinderen er bv. zo kunnen uitzien:



2-point crossover met wisselposities 3 en 6:



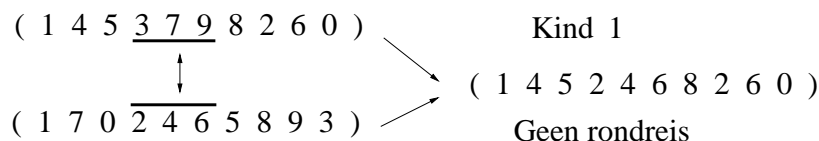
**Belangrijk:** De net geziene principes voor mutatie (lettertekens uit het alfabet vervangen) en crossover (vervangen van 2 deelstrings op dezelfde positie van de twee ouders) zijn handleidingen hoe je het in de meeste gevallen kan doen, maar soms moet je de ideeën die erachter staan op een iets andere manier implementeren.

Belangrijk voor mutatie is *een individu maar een beetje te wijzigen* en voor crossover is het vooral dat *karakteristieken* van beide ouders in het kind aanwezig zijn. Een voorbeeld waar de standaardtechnieken niet goed toegepast kunnen worden:

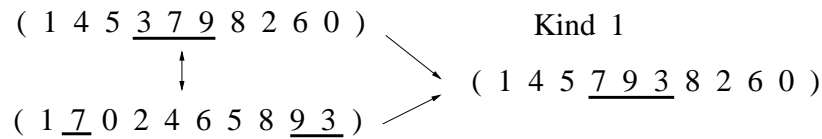
**Voorbeeld 7** Wij zoeken een goedkoopste rondreis in een gewogen complete graaf. Wij coderen elke rondreis door een string met de nummers van de toppen.

Als wij een teken voor een stad door een ander teken vervangen, hebben wij één stad twee keer en een andere niet meer – dat werkt dus niet. Hier zouden wij mutatie bv. kunnen implementeren door toevallig een positie te kiezen en de stad met de volgende van volgorde te verwisselen.

Als wij crossover op de gewone manier doen, zal het resultaat vaak geen rondreis zijn:



Een andere manier een crossover te definiëren is als volgt: Je kiest nog altijd een deelstring van de eerste ouder, maar vervangt hem niet door de deelstring op dezelfde posities van de andere ouder, maar herordent de deelstring zo dat de steden in deze deelstring in de volgorde van de tweede ouder worden bezocht. Op deze manier kan je garanderen dat elk kind ook tot een punt (in dit geval een rondreis) behoort.



### Selectie:

Terug naar het algemeen geval: Wij hebben twee stappen waar selectie gebeurt: één keer wanneer het algoritme moet kiezen welk individu kinderen verwekt en één keer waar het beslist wie deel van de volgende populatie is (dus welk individu moet sterven). Zoals voor alles bestaan ook voor deze selectiestappen meerdere mogelijkheden. Een gretige aanzet waar bijvoorbeeld altijd alleen maar de besten (die met de grootste fitness of met andere woorden *de fitste*) kinderen verwekken en overleven, bleken vaak niet goed te werken. Het is beter deze beslissingen ook gerandomiseerd te nemen: Hoe groter de *fitness* (waarde van de doelfunctie) van een individu is, hoe groter de kans dat het kinderen verwekt en hoe groter de kans dat het ook deel uitmaakt van de volgende populatie. Maar het is belangrijk dat ook individuen met een kleine fitness kans maken lid van de volgende populatie te zijn en kinderen te verwekken.

Ook individuen met een relatief lage fitness kunnen belangrijke eigenschappen bevatten die ze kunnen bijdragen aan de populatie en die na een mutatie of door crossover met andere individuen tot een bijzonder goede oplossing kunnen leiden. Aan de andere kant zijn meerdere heel goede – maar ook heel gelijkaardige – individuen niet nuttig, omdat ze alleen maar dezelfde eigenschappen kunnen bijdragen. Genetische algoritmen werken alleen goed als je voldoende diversiteit in jouw populatie hebt, dus zoiets als een rijke genpool.

Als wij al deze delen die wij net hebben besproken samenvatten dan hebben wij het volgende algoritme:

### Algoritme 7 genetisch algoritme

*Gegeven een verzameling  $M$ , een fitnessfunctie  $f()$  waarvan wij het maximum zoeken.*

Wij hebben een codering  $c : M \rightarrow A^*$  voor een alfabet  $A$ , de parameters en methoden voor crossover, mutatie en selectie vastgelegd en beslist hoe groot de populatie moet zijn.

(a) kies een beginpopulatie (bv. op een toevallige manier).

(b) Herhaal de volgende stappen tot aan de stopvoorwaarde is voldaan.

- selecteer individuen om kinderen te verwekken
- vorm paren en verwek kinderen met hulp van crossover
- pas mutatie toe op de kinderen
- voeg de nieuwe kinderen toe tot de populatie
- selecteer de nieuwe populatie (dus beslis wie sterft zodat de populatie niet groeit).

(c) Het resultaat is het beste individu uit de laatste populatie.

**Oefening 23** Ontwerp een genetisch algoritme voor het probleem dat je een deelgraaf met minimaal gewicht van een gewogen complete graaf zoekt (waar de gewichten ook negatief kunnen zijn) zodat elke top die in de deelgraaf zit graad 2 in de deelgraaf heeft.

Het handelsreizigersprobleem is dus een speciaal geval van dit probleem – het geval waar de deelgraaf opspannend en samenhangend moet zijn (en de gewichten positief zijn, maar dat maakt niet echt een verschil).

Geef alle definities en parameters die nodig zijn om het algoritme volledig te beschrijven maar het is niet nodig dat ze ook zo gekozen zijn dat het algoritme goed werkt. Daarvoor zou je het echt moeten implementeren en testen laten draaien. Het algoritme moet alleen maar in principe werken en het optimum kunnen vinden (ook al zou dat veel te lang duren en te veel pogingen vragen). De bedoeling is gewoon om aan te tonen dat het principe verstaan is.

**Oefening 24** Ontwerp een genetisch algoritme voor het probleem een maximale matching te vinden. Gebruik als  $M$  de verzameling van alle matchings. Hoe kan je nu  $c()$ , mutatie en crossover definiëren?

**Oefening 25** Beschrijf een genetisch algoritme voor het volgende probleem. Het is voldoende alle parameters en constructies te beschrijven en het basisalgoritme te geven.

**Probleem:** In een magazijn staan  $n$  goederen  $g_1, \dots, g_n$  te koop. Voor elk goed  $g_i$  is de prijs  $p(g_i)$  die je zelf moet betalen gekend en ook de prijs  $e(g_i)$  die je kan krijgen als je het aan een eindverbruiker doorverkoopt.



*Het doel is nu voor een gegeven getal  $M$  een verzameling van goederen te vinden zodat de som van de prijzen ten hoogste  $M$  is en jouw winst optimaal.*

## **2.6 Afsluitende opmerkingen**

Hoewel wij er meerdere lessen aan hebben besteed, waren het alleen maar de oppervlakkige ideeën van de metaheuristieken die wij konden bespreken. En bovendien zijn er nog veel metaheuristieken waarvan wij zelfs de ideeën niet hebben besproken (ant colony optimization, particle swarm optimization, threshold accepting, tabu search, great deluge algorithm, asynchronous teams, de klasse van memetic algoritmen (genetische algoritmen zoals gezien zijn een voorbeeld daarvan) etc.) Over elk van de metaheuristieken die wij hebben besproken zijn er meer artikels gepubliceerd dan jullie ooit kunnen (of willen) lezen.

Als jullie een probleem moeten oplossen waarvoor een metaheuristiek een goede keuze is, dan is het vooral belangrijk het probleem goed te verstaan en de functies (bv.  $N()$ ) en parameters op een manier te kiezen die ideaal is voor dit **specifieke** probleem. Vooral in het geval van genetische algoritmen zijn er ongelooflijk veel stappen die op verschillende manieren gedaan kunnen worden en parameters die gekozen moeten worden (het best door veel tests te draaien). Hoewel ook voor specifieke problemen heel goede en mooie ideeën ontwikkeld werden om een optimale prestatie van de metaheuristieken te garanderen, zijn deze ideeën vaak te specifiek om voor een les interessant te zijn.

Wat algemeen geldt, is dat de onderdelen van de metaheuristieken heel snel moeten zijn. Dat geldt natuurlijk vooral voor onderdelen van het algoritme die heel vaak uitgevoerd moeten worden (bv. de berekening van buuren). Daardoor zijn **te** ingewikkelde definities van bv. de buurfuncties een slechte keuze.

In sommige gevallen waar je heel goed weet waar *ongeveer* de goede oplossingen zijn, is het beter een startpunt te kiezen in plaats van een toevallig startpunt te nemen.

### 3 Algoritmen voor slim gebruik van het geheugen

Een tijdje geleden was er in DA2 een project waar verschillende *AVL-achtige* bomen geïmplementeerd moesten worden. De datastructuur was identiek en de routines die sleutels zochten waren ook dezelfde. Volgens onze ideeën over tijdsverbruik zou dus het aantal bezochte toppen voor een reeks van zoekoperaties een goede manier zijn om de tijdscomplexiteit te meten. Maar de resultaten waren als volgt (2.000.000 sleutels in de boom en 20.000.000 keer een toevallige sleutel opzoeken):

algoritme	tijd voor opbouwen	aantal bezochte toppen (zoeken)	tijd voor zoeken
1	3.75	426.500.000	33.3
2	8.07	420.300.000	37.6

Hoewel de boom dus beter gebalanceerd was en dezelfde routines de opzoekingen deden, vroeg het meer tijd om de sleutels op te zoeken – er klopte dus iets niet met ons model. Een idee was dat het iets met de cache te maken had. De slechter gebalanceerde boom werd minder herbalanceerd en misschien was de manier waarop het geheugen gealloceerd werd beter omdat toppen die dicht bij elkaar in de boom zitten ook dicht bij elkaar in het geheugen zitten en dus samen in de cache plaats kunnen vinden.

Om dat te testen werden de toppen van de bomen na het toevoegen in het geheugen verplaatst zodat ten minste de toppen die dicht bij de wortel waren allemaal ook dicht bij elkaar zaten. De structuur van de bomen was nog dezelfde (de bomen waren isomorf aan de bomen voor het *verhuizen*). Het resultaat was dat beide algoritmen (nog steeds met dezelfde routines) sneller waren – maar deze keer was de verhouding (ongeveer) zoals verwacht:

algoritme	aantal bezochte toppen (zoeken)	tijd voor zoeken
1	426.500.000	27.7
2	420.300.000	26.8

De modellen die wij gebruikt hebben om te bepalen of een algoritme efficiënt is of niet zijn dus in dit geval niet echt nauwkeurig – en hier gaat het alleen maar om het verschil tussen de cache en het geheugen van de computer. Zodra er zoveel bestanden zijn of de records zo groot zijn dat de boom op de harde schijf wordt bijgehouden, zijn het de leesoperaties op de harde schijf die de snelheid bepalen en niet het aantal vergelijkingen.

De 2-3-bomen die wij in DA2 hebben gezien, gaan al een stukje in deze richting. Als wij het aantal toppen als het aantal leesoperaties zien dan

hebben 2-3-bomen een voordeel in vergelijking met alle binaire zoekbomen omdat de diepte het aantal leesoperaties en daardoor de snelheid bepaalt. Door de grotere vertakking hebben 2-3-bomen een kleinere diepte dan binaire zoekbomen (behalve in het ene geval waar een 2-3-boom een complete binaire boom kan zijn) en dat gaat nu ook onze strategie zijn: de vertakking zo groot mogelijk maken zodat zo weinig mogelijk leesoperaties nodig zijn. Daarbij gaat het niet alleen om zoekbomen – je hebt ook een vertakking als het bv. om recursie gaat. Ons eerste voorbeeld gaat er één zijn waar de vertakking de recursie beschrijft en niet de datastructuur.

### 3.1 Hashing

Hashing is een efficiënte manier om sleutels op te slaan en op te zoeken. Het probleem is alleen dat er botsingen kunnen zijn – dan kan het gebeuren dat hashing niet meer efficiënt is. Je kan botsingen voorkomen door een hash-tabel te kiezen die voldoende groot is om de kans op botsingen klein te houden. Het probleem is dat je natuurlijk niet te veel ruimte wil verspillen en dat je voor de geziene manieren om met hashtabellen te werken op voorhand moet vastleggen hoe groot de tabel is.

Het zou dus mooi zijn als je de tabel zou kunnen uitbreiden tijdens het gebruik als je ziet dat er te weinig ruimte is. Het onderwerp van dit deel zijn twee manieren om met hashtabellen te werken die dat op een efficiënte manier toelaten.

**Oefening 26** *Stel dat je closed hashing wil toepassen. Om geen probleem met te veel botsingen te hebben, wordt elke keer dat de laadfactor boven de 70% ligt de hashtabel uitgebreid. Wij hebben al in DA2 gezien dat het uitbreiden van een array niet efficiënt kan gebeuren als je elke keer een vast aantal vakken toevoegt maar dat het nodig is de grootte met een getal (in de meeste gevallen 2) te vermenigvuldigen. Stel dat je hier (omdat je toch al ruimte verspilt) als factor maar 1.5 kiest. Uitbreiden betekent hier wel dat de tabel niet alleen gekopieerd moet worden, maar ook dat voor alle sleutels de hashfunctie opnieuw berekend moet worden.*

*Stel voor het volgende dat er geen botsingen zijn:*

- *Hoe duur is een toevoegbewerking in het slechtste geval?*
- *Hoe duur is een reeks van  $n$  toevoegbewerking op een initieel lege hash-tabel met 100 elementen in het slechtste geval? (Herinnering: dat hebben wij ook de geamortiseerde kost van deze reeks van bewerkingen genoemd.)*

In dit deel zijn we er vooral in geïnteresseerd als de records die je wil opslaan heel groot zijn – zo groot dat je ze niet in het geheugen kan houden maar dat je ze op de harde schijf moet plaatsen. Over hoe dat precies gebeurt gaan wij het niet hebben – dat laten wij gewoon over aan de controller van de harde schijf en aan het besturingssysteem. Bovendien hebben jullie daar al iets over gezien. Wij willen het aantal lees/schrijf-operaties op de harde schijf beperken. Een bestand op de harde schijf kan je – anders dan een array werkt in het geheugen – langer maken zonder het te moeten kopiëren (dat zou heel veel leesoperaties vragen). Wij werken hier gewoon met een model waarvan wij stellen dat je een geïndexeerde lijst van *vakjes* hebt zodat je elk vakje in één stap kan lezen en in één stap er iets in schrijven. Bovendien stellen wij dat je de lijst langer kan maken en dat de kosten om er een vakje aan toe te voegen dezelfde zijn als voor een gewone schrijfoperatie. Hoe het besturingssysteem ervoor zorgt dat dat kan, doet er **in dit geval** niet toe. Een model is nuttig als je gebaseerd op dit model snelle algoritmen ontwikkeld, waarbij *snel* hier de echte tijd in de praktijk betekent – dat is altijd ons doel. Ook al zijn in ons model sommige dingen zeker te gemakkelijk voorgesteld, blijkt het toch een in de praktijk nuttig model.

### 3.1.1 Extendible hashing

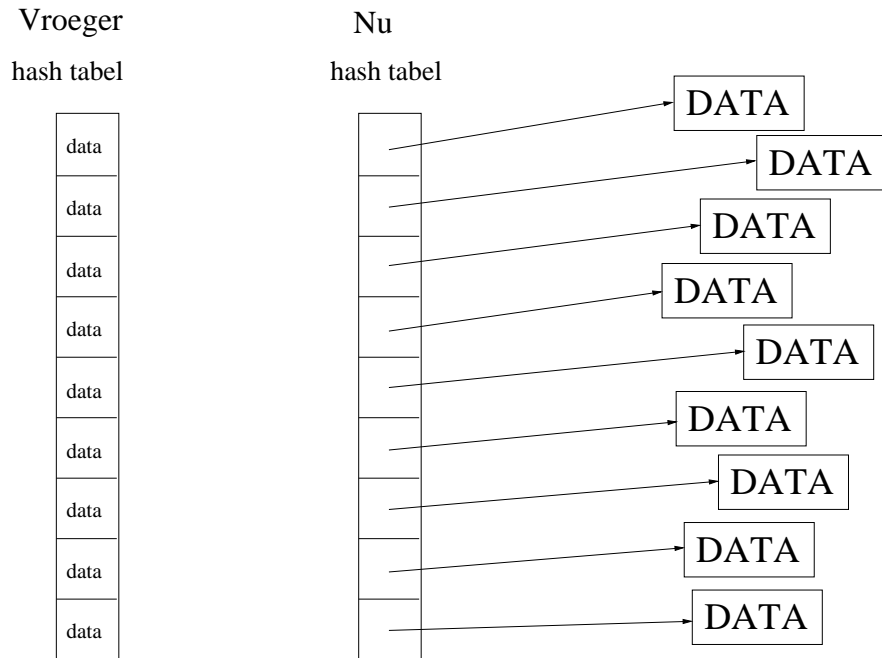
Extendible hashing werkt met een model dat een beetje verschilt van het gewone model waar je de records (de data) direct in een tabel opslaat. Wij werken met een tabel van pointers die naar de echte data wijzen. In een reële toepassing kan deze *echte data* natuurlijk ook maar een deel van de echte data zijn of alleen een identificatienummer dat nog een pointer naar de hele record bevat – dus de echte data alleen maar representeert.

Als je de records in het geheugen kan houden, kan deze manier van werken ook voordelen hebben als je niet met extendible hashing werkt: het zou kunnen dat je moeilijkheden hebt een samenhangend geheugenblok te krijgen om alle records te indexeren. Als je met pointers werkt, is dat natuurlijk geen probleem. Bovendien zou het ook efficiënter zijn de array groter te maken: Wij zouden alleen de pointers moeten kopiëren en niet de (veel grotere) records. Maar zonder verdere wijzigingen zouden wij de data nog altijd moeten rehashen of tenminste lezen en dat is natuurlijk ook vrij duur – en als het over de harde schijf gaat onaanvaardbaar.

Maar nu gaan wij het over het geval hebben waar de records op de harde schijf zitten.

De array met de pointers kan zeker nog in het geheugen gehouden worden, dus hebben wij daarvoor geen leesoperaties op de schijf nodig.

Wij werken met buckets (emmers) waarbij in elke bucket een constant en



Figuur 1: Een hash tabel en een hash tabel met pointers naar buckets.

vast aantal records terecht kan komen – bv. 5. De reden is dat de emmers blokken voorstellen die je in één stap kan lezen. Het feit dat je dan misschien de gezochte record nog moet zoeken als er meerdere records per emmer zijn, doet er in ons model niet toe: dat gebeurt in het geheugen en kan dus efficiënt gedaan worden. Wij werken met een hash-functie  $h()$  waarvan wij stellen dat de waarden tussen 0 en  $2^n - 1$  voor een zekere  $n$  liggen en dat ze zeker voldoende groot zijn om goed te kunnen hashen – bv. waarden tussen 0 en  $2^{32} - 1$ . Bovendien stellen wij dat de hashfunctie altijd met hetzelfde aantal  $n$  bits voorgesteld wordt – dus misschien ook met nullen in het begin. Wij zullen in onze voorbeelden altijd de hash-waarden in de emmers opslaan. In de realiteit zijn dat natuurlijk de records en niet de hashwaarden, maar zo kan je zien wat er precies gebeurt terwijl artificiële data die wij voor de voorbeelden zouden gebruiken niet echt zou bijdragen tot het verstaan van de principes. Je zou in de realiteit bovendien **ook** nog de hashwaarden kunnen opslaan om die niet te moeten herberekenen, maar omdat nooit veel data opnieuw gehasht moet worden en omdat wij alleen de lees- en schrijf-bewerkingen willen minimaliseren, maakt dat geen (groot) verschil.

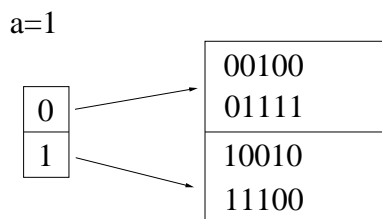
Arrayindices schrijven wij ook vaak in binaire vorm omdat de methodes de binaire voorstelling van de indices gebruiken. Natuurlijk zijn dat ook ge-

wone getallen – alleen dat ze binair voorgesteld worden om gemakkelijker te verstaan in welk vakje een record terecht moet komen.

- De grootte van onze pointerarray is altijd een macht van 2 – wij gebruiken de variabele  $a$  om dat bij te houden, de grootte is  $2^a$ .

Wij evalueren altijd een zeker aantal – dezelfde  $a$  – bits van onze hashfunctie. Inderdaad gebruiken wij dus niet  $h()$  om de index van het juiste vakje in onze tabel te vinden maar  $(h())|_a$ . Wij schrijven altijd  $x|_a$  voor het getal dat door de eerste  $a$  bits in de binaire voorstelling van  $x$  voorgesteld wordt, waarbij wij misschien met nullen in het begin werken om ervoor te zorgen dat de voorstelling van elk getal hetzelfde aantal bits heeft. Als  $x = 010011$  in binaire voorstelling (dus 19) dan is  $x|_3 = 010$  dus gelijk aan 2. Als wij  $a$  bits evalueren, hebben wij  $2^a$  mogelijke waarden  $0, \dots, 2^a - 1$  die de indices van onze (pointer) hash-tabel vormen. *In principe* kunnen wij dus pointers naar  $2^a$  buckets bijhouden.

Een voorbeeld met  $a = 1$  zien jullie in Figuur 2



Figuur 2: Extendible hashing met  $a = 1$  en twee records per emmer.

Maar inderdaad zullen wij niet alle buckets wijzigen als de pointerarray groter gemaakt wordt. Dat is **heel** belangrijk omdat het veel dure leesoperaties zou vragen en de uitbreidingsstap onaanvaardbaar duur zou maken. Dat betekent dat wij voor sommige buckets niet  $a$  bits moeten evalueren maar minder. Wij houden dus voor elke bucket  $b$  een getal  $a_b$  bij dat zegt hoeveel bits voor deze bucket werden geëvalueerd toen hij aangemaakt werd. Het is duidelijk dat als er meer pointers zijn dan buckets dat sommige pointers naar dezelfde bucket moeten wijzen. Figuur 4 toont een situatie waar er één bucket al werd uitgebreid en één niet. Inderdaad zullen er altijd  $2^{a-a_b}$  pointers naar een bucket  $b$  wijzen.

Maar nu dat wij al ongeveer weten hoe de situatie er kan uitzien, moeten wij precies beschrijven hoe extendible hashing werkt:

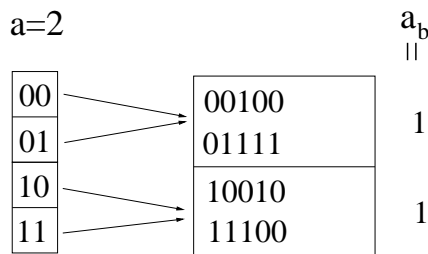
Je begint met een vaste  $a$ . In onze voorbeelden beginnen wij met  $a = 1$  omdat het anders vrij moeilijk is om te tekenen, maar in de realiteit zal  $a$

zeker ten minste 10 zijn. Wij hebben dan een pointerarray met  $2^a$  pointers die indices  $0, \dots, 2^a - 1$  hebben en naar  $2^a$  buckets wijzen waarin wij  $g$  records kunnen plaatsen. Voor de voorbeelden kiezen wij  $g = 2$ . Voor elke emmer  $b$  is in het begin  $a_b = a$ . In het begin – waar voor elke emmer geldt dat  $a = a_b$  – wijzen dus inderdaad  $2^{a-a_b} = 1$  pointers naar elke emmer.

- Toevoegen en opzoeken gebeurt op de volgende manier: je berekent de hashwaarde  $h(r)$  van de record  $r$  die je wil toevoegen of opzoeken en kijkt in de tabel waar vakje  $(h(r))|_a$  naartoe wijst. Daar moet je record  $r$  plaatsen of opzoeken. Opzoeken is nooit een probleem maar als je de record wil toevoegen en de emmer is vol is er wel een probleem...

Als de emmer  $b$  waar je iets wil toevoegen volzit, willen wij een nieuwe emmer gebruiken en de inhoud verdelen. Maar voor 2 emmers hebben wij ook 2 pointers nodig. Er zijn twee mogelijkheden:

$a_b = a$ : In dit geval is er maar één pointer die naar deze emmer wijst. Dat is het slechtste geval. Wij moeten dus eerst onze pointerarray uitbreiden om ervoor te zorgen dat ten minste 2 pointers naar deze emmer wijzen. Wij verdubbelen de pointerarray. De variabele  $a$  is nu 1 groter (dus  $a = a^{oud} + 1$ ). Als de nieuwe pointerarray  $p_n[]$  is en de oude  $p_o[]$  dan vullen wij  $p_n[]$  met dezelfde pointers als  $p_o[]$ : voor  $0 \leq i < 2^a$  geldt  $p_n[i] = p_o[\lfloor i/2 \rfloor]$ . Het resultaat voor de eerste keer verdubbelen zien jullie in Figuur 3. Nu wijzen dus twee keer zo veel pointers naar elke emmer als voor het verdubbelen.



Figuur 3: De pointerarray werd verdubbeld en werkt met  $a = 2$  maar de emmers werken nog met  $a_b = 1$ .

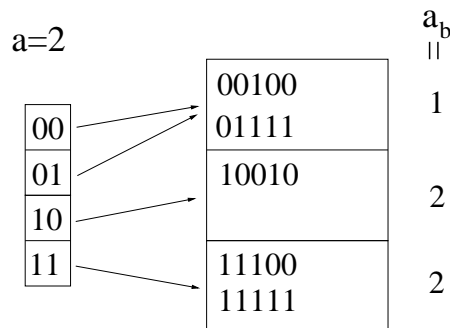
Wij hebben er dus voor gezorgd dat wij altijd in de volgende situatie terechtkomen:

$a_b < a$ : Dus wijzen nu  $2^{a-a_b} \geq 2$  pointers naar de volle emmer  $b$  waar wij de nieuwe record willen plaatsen. De eerste  $a_b$  bits van de hashwaarde van

elke record zijn gelijk – stel dat het  $b_1 \dots b_{a_b}$  zijn. Wij maken nu een nieuwe emmer  $b'$  aan en verdelen de records uit  $b$  zo dat alle records  $r$  met bit nummer  $a_b + 1$  gelijk aan 0 (of  $2 \cdot (h(r)|_{a_b}) = h(r)|_{a_b+1}$ ) in de oude emmer  $b$  terechtkomen. De eerste  $a_b + 1$  bits van de hashwaarden van deze records zijn dus “ $b_1 \dots b_{a_b} 0$ ”. De anderen (bit nummer  $a_b + 1$  gelijk aan 1, dus  $2 \cdot (h(r)|_{a_b}) + 1 = h(r)|_{a_b+1}$  en de eerste  $a_b + 1$  bits van de hashwaarde zijn “ $b_1 \dots b_{a_b} 1$ ”) plaatsen wij in de nieuwe emmer  $b'$ . Dan plaatsen wij de nieuwe record volgens dezelfde regel (als dat kan. . .). De pointers in de pointerarray met een index  $i$  waarvoor geldt dat  $i|_{a_b+1} = 2 \cdot (h(r)|_{a_b}) + 1$  (de eerste  $a_b + 1$  bits van de index zijn “ $b_1 \dots b_{a_b} 1$ ”) moeten nu gewijzigd worden zodat ze naar  $b'$  wijzen. De waarden van  $a_b$  en  $a_{b'}$  worden de oude waarde van  $a_b$  plus 1.

**Oefening 27** Bereken uit  $a$ ,  $a_b$  en  $h(r)$  de indices van de pointers die gewijzigd moeten worden. Geef een formule.

In Figuur 4 zien jullie het resultaat als wij op deze manier 11111 toevoegen: de emmer waarop pointer 11 wees moest herverdeeld worden en omdat  $h(r)|_{a_b} = 1$  (alleen het eerste bit wordt geevalueerd van 11111) wordt pointer  $2 \cdot 1 + 1 = 3$  gewijzigd en wijst naar de nieuwe emmer.

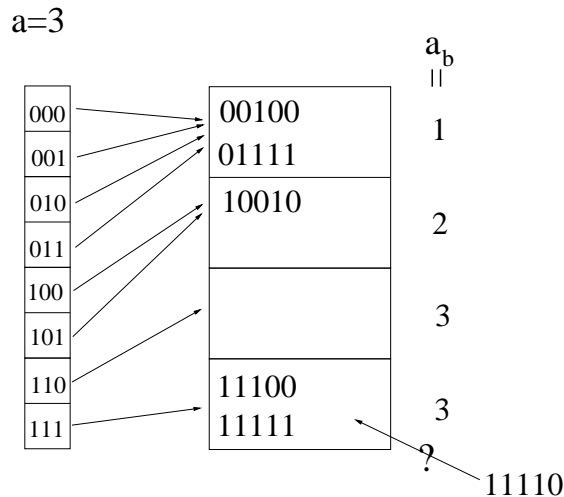


Figuur 4: Extendible hashing met  $a = 2$  maar één emmer werkt nog met  $a_b = 1$ .

Maar inderdaad kan er een probleem zijn. Wij hebben de oude records in de emmer en de nieuwe record herverdeeld over twee emmers. Maar als wij pech hebben, kan het gebeuren dat bij het herverdelen één emmer leeg blijft en de andere dan nog altijd geen vrije ruimte voor de nieuwe record heeft. Dat zou bv. gebeuren als je nu sleutel 11110 toevoegde. Figuur 5 toont de datastructuur na het dubbel en herverdelen van de volle emmer.

In dit geval moeten wij nog eens verdubbelen en hopen dat het probleem dan opgelost is. Hashing kan natuurlijk altijd problemen veroorzaken als





Figuur 5: Extendible hashing met  $a = 3$  na het verdubbelen maar toch kan een record nog niet toegevoegd worden.

door een slechte hash-functie – of gewoon heel veel pech – de waarden slecht verdeeld zijn. Dit zou ook hier de oorzaak zijn. Met een goede hash-functie en niet te weinig records per emmer zal de load in de praktijk gelijkmatig verdeeld zijn en dus de pointerarray niet overbodig groot zijn. Als in één emmer  $x$  records terecht kunnen komen en de hash-functie kent meer dan  $x$  keer dezelfde waarde toe, zou het uitbreiden oneindig doorgaan (of precies: het zou doorgaan tot er geen bits meer zijn om te evalueren en zou dan vaststellen dat de methode niet werkt). Een goede hashfunctie is dus heel belangrijk.

**Het** belangrijke voordeel qua efficiëntie is dat je alleen maar naar de emmer moet kijken die te vol zit. Dat zorgt ervoor dat je misschien met maar één leesoperatie en twee schrijfoperaties kan werken als de grootte van de emmers zo is dat hij in één stap gelezen kan worden. Naar de andere emmers moet je zelfs niet kijken. Het uitbreiden van de pointerarray gebeurt in het geheugen – dat tellen wij dus niet mee als het om lees/schrijf-operaties gaat. Bovendien moet de array natuurlijk **veel** minder vaak gedubbeld worden dan emmers herverdeeld moeten worden – als er  $m$  keer een emmer herverdeeld moet worden dan is het aantal dubbeloperaties bij een gelijkmatige verdeling van de hashwaarden van de orde  $\log m$ . Maar het dubbelen kan ook op een (tenminste geamortiseerd) efficiënte manier gebeuren. Of precies:

**Oefening 28** *Deze oefening is niet gemakkelijk (vraag de assistent je tips te geven als je niet weet hoe je moet beginnen) maar het is wel een nuttige*

oefening om een beetje DA2 te herhalen en vooral om zichzelf duidelijk te maken hoe extendible hashing werkt:

Het verdubbelen van de pointerarray lijkt natuurlijk sterk op het verdubbelen van een array, dat jullie in DA2 hebben gezien. Maar er zijn toch verschillen: hier wordt de array direct volgeschreven en bovendien worden de records niet gekopieerd en de samenhang tussen het aantal records en de grootte van de pointerarray is ook niet zo duidelijk.

- Wat is de geamortiseerde kost van een reeks van  $n$  bewerkingen op een initiëel lege extendible hashing tabel als je geen voorwaarden aan de hashfunctie oplegt (behalve natuurlijk dat ze oneindig veel bits genereert – anders werkt het niet voor arbitrair grote  $n$ )?
- Gegeven  $q > 0$ . Wat is de geamortiseerde kost van een reeks van  $n$  bewerkingen op een initiëel lege extendible hashing tabel als gegarandeerd is dat altijd als er een tabel met  $p$  pointers uitgebreid wordt er ook ten minste  $p \cdot q$  records zijn? De  $O()$  notatie is voldoende maar een bewijs is (natuurlijk) vereist. (Inderdaad zou het voldoende zijn dat voor de datastructuur na  $n$  stappen te eisen.)

**Oefening 29** Stel dat je met extendible hashing ook identieke records – dus zeker met gelijke sleutels – wilt opslaan. Wanneer kan dat en wanneer is er een probleem? Hoe kan je het probleem oplossen?

**Oefening 30** Voeg de sleutels 3, 7, 12, 45, 44 en 1 toe aan een lege extendible hashing tabel met twee emmers in het begin. Er kunnen 2 sleutels (records) per emmer geplaatst worden.

Wij gebruiken een voorstelling van de hashcodes van  $h()$  als binaire getallen met 5 bits en wij hebben

$h(3) = 0$  dus binair 00000

$h(7) = 31$  dus binair 11111

$h(12) = 9$  dus binair 01001

$h(45) = 17$  dus binair 10001

$h(44) = 2$  dus binair 00010

$h(1) = 1$  dus binair 00001

### 3.1.2 Linear hashing

In het geval van extendible hashing kan een enkele uitbreidingsstap *in principe* wel relatief duur zijn – ook al is de geamortiseerde kost (met een goede hashfunctie) heel goed. Het effect van een slechte hashfunctie – of heel veel

pech – is dat meerdere keren uitgebreid moet worden of in het ergste geval zelfs de grens van de uitbreidingsmogelijkheden bereikt wordt. Ook het verdubbelen van de pointerarray kan je voor zeer grote arrays niet meer verwaarlozen. Linear hashing gebruikt uitbreidingsstappen waar het slechtste geval goedkoper is. De klemtoon ligt hier op de laadfactor en niet op een te volle emmer. Je moet afhankelijk van de toepassing beslissen of extendible hashing of linear hashing (of gewoon hashing) beter is.

Hier gebruiken wij het voordeel van ons model van de harde schijf dat de tabel zonder te kopiëren uitgebreid kan worden. Het is dus mogelijk een enkele emmer toe te voegen zonder de andere emmers te moeten kopiëren.

De belangrijkste bedoeling is ook hier tijdens het uitbreiden het aantal lees- en schrijf-operaties op de harde schijf te minimaliseren – dus naar zo weinig mogelijk emmers te moeten kijken.

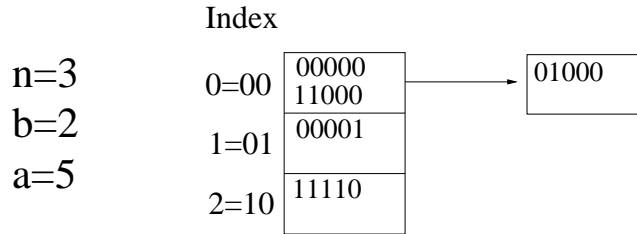
Ook hier zullen wij in de voorbeelden de hash-keys in de emmers plaatsen omdat dat gemakkelijker is om te verstaan. In een toepassing zitten daar natuurlijk de records.

- Je houdt altijd bij hoeveel emmers je hebt (dat noemen wij altijd  $n$ ) en hoeveel elementen in de array zitten (dat noemen wij  $a$ ). Het aantal records dat je in één emmer kan plaatsen noemen wij  $g$  (voor grootte).

Een groot verschil met extendible hashing is dat wij er hier niet op letten of een emmer te vol zit of niet. Als iets niet meer geplaatst kan worden, werken wij gewoon met overstroommemmers. Natuurlijk maakt dat hashen misschien minder efficiënt – vooral als voor de overstroommemmers een extra leesoperatie nodig is – maar als de hashfunctie goed en de laadfactor  $a/(n \cdot g)$  niet te groot is, zullen er normaal niet veel overstroom-emmers zijn. Als wij schrijven dat in een *emmer* gezocht moet worden, bedoelen wij inderdaad de hele emmerketting – dus de emmer samen met zijn overstroommemmers. Ons doel is dus de laadfactor relatief klein te houden om efficiënt hashen te kunnen garanderen.

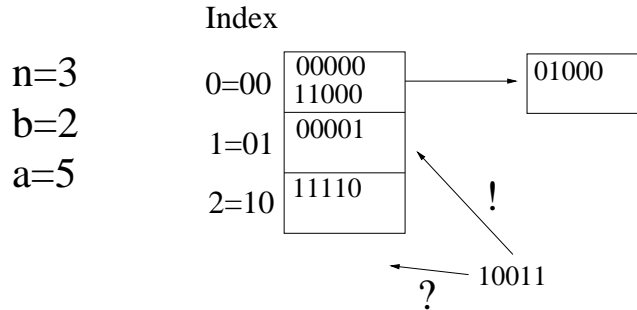
Wij gebruiken opnieuw de binaire voorstelling van de hash-waarden. Als er  $n$  emmers zijn, hebben wij  $b(n) = \lceil \log n \rceil$  bits nodig om de indices van de emmers voor te stellen. Hier gebruiken wij de **laatste**  $b$  bits van de hashfunctie  $h()$  en schrijven  $x|b$  als wij de laatste  $b$  bits van de binaire voorstelling van een getal  $x$  willen gebruiken. Wij stellen altijd dat  $n \geq 2$  en dus  $b \geq 1$ . Wij proberen een record  $r$  in emmer  $h(r)|b$  te plaatsen. Het enige probleem dat zou kunnen opduiken is dat  $h(r)|b \geq n$ . Alleen als  $n$  een macht van twee is, is  $b(n) = \log n$ . Dan kan dat niet gebeuren – anders wel!

Als je in ons voorbeeld een record met hashwaarde 10011 wil plaatsen, neem je de laatste  $b = 2$  bits, dus 11 – het record moet dus in vakje 3 geplaatst



Figuur 6: Een kleine linear hashing tabel.

worden – maar dat is er niet. In zulke gevallen plaatsen wij het record in vakje  $h(r)|^b \bmod 2^{b-1} = h(r)|^{b-1}$ .



Figuur 7: Het plaatsen van een element waar de eerste keuze voor de index te groot is.

De functie  $h_n(r)$  die de index bepaalt waar het element  $r$  geplaatst moet worden is dus gegeven door (merk op dat  $b$  gewoon een functie van  $n$  is):

$$h_n(r) = \begin{cases} h(r)|^b & \text{als } h(r)|^b < n \\ h(r)|^{b-1} & \text{als } h(r)|^b \geq n \end{cases}$$

Op deze manier hebben wij een manier van hashing die ook zonder uitbreiden kan werken en door de overstromemmers ook met een grote laadfactor nog werkt – hoewel de efficiëntie dan slechter wordt. Het lineaire zoeken in de overstromemmers zou misschien veel leesoperaties vragen.

Maar nu gaan wij beschrijven hoe wij ervoor zorgen dat de laadfactor  $a/(n \cdot g)$  nooit groter wordt dan een bovengrens  $L$  die wij op voorhand vastleggen. Het getal  $n \cdot g$  is het grootste aantal records dat je kan plaatsen zonder overstromemmers. Hoe dichter  $a$  dus bij  $n \cdot g$  ligt hoe groter de kans dat je overstromemmers moet gebruiken. Als  $a > n \cdot g$  is het zelfs zeker dat je overstromemmers nodig hebt! Wij zullen dus normaal  $L < 1$  kiezen – bv.  $L = 0.7$ .

Stel nu dat wij een nieuw element toevoegen en daardoor de laadfactor te groot wordt, dus  $a/(n \cdot g) > L$ . Dan moeten wij onze tabel uitbreiden. Daardoor wordt  $n$  groter en de laadfactor kleiner. Wij nemen er altijd maar één nieuw element bij.

Maar dat betekent natuurlijk dat onze hashfunctie  $h_n(r)$  verandert! Wij moeten dus bepalen welke records  $r$  nu op een foute plaats zitten dus waarvoor geldt  $h_{n_{\text{nieuw}}}(r) \neq h_{n_{\text{oud}}}(r)$  of als  $n$  voor de oude  $n$  staat:  $h_{n+1}(r) \neq h_n(r)$ . Er zijn twee gevallen: ofwel  $b$  verandert ook – ofwel niet.

$b(n) = b(n+1)$ : Wij schrijven gewoon  $b$  voor  $b(n) = b(n+1)$ . Als je naar de definitie van  $h_n(r)$  kijkt, kan een verschil er alleen door veroorzaakt zijn dat verschillende gevallen van de definitie worden toegepast. Dat gebeurt precies voor  $h|b(r) = n = n_{\text{oud}}$ . Vroeger was de waarde  $n|^{b-1}$  en nu is hij  $n$ . De enige records die op een foute plaats kunnen zitten, zijn dus de records in emmer  $n|^{b-1}$ . Die moeten herverdeeld worden.

$b(n) = b(n+1) - 1$ : Dus  $n = n_{\text{oud}} = 2^{b_{\text{oud}}}$  Wij schrijven  $b$  voor  $b_{\text{oud}}$ . Nu worden  $b+1$  bits van  $h(r)$  geëvalueerd (dus  $h(r)|^{b+1}$  gebruikt). Als de eerste bit 0 is, geldt  $h(r)|^{b+1} = h(r)|^b < n < n+1$ . Records met deze hashwaarde werden dus vroeger en nu op dezelfde manier geplaatst. Maar als de eerste bit (die voor  $2^b$  staat) 1 is, geldt  $h(r)|^{b+1} \geq 2^b = n$ . Als  $h(r)|^{b+1} > 2^b = n$  dan geldt  $h(r)|^{b+1} \geq n+1$  – dus is  $h_{n+1}(r) = h(r)|^b = h_n(r)$  (omdat in dit geval altijd geldt dat  $h(r)|^b < n$  – zelfs als alle bits 1 zijn.). Records met deze hashwaarde moeten dus niet herverdeeld worden. Het blijft alleen maar  $h(r)|^{b+1} = 2^b = n$  dus  $h(r)|^b = 0$  om herverdeeld te worden.

Je hoeft in de twee gevallen dus geen verschillende formules te gebruiken – in elk geval (behalve het uitzonderlijke geval  $b_{\text{oud}} = 0$  dat je voor echte toepassingen natuurlijk nooit hebt) is de emmer die herverdeeld moet worden emmer  $n_{\text{oud}}|^{b_{\text{oud}}-1}$ . De twee gevallen die wij besproken hebben, zijn gewoon bewijzen dat deze formule in beide gevallen geldt.

De emmer die herverdeeld wordt, moet niet bijzonder vol zitten – hij kan (in principe) zelfs helemaal leeg zijn!

**Oefening 31** *Als je veronderstelt dat inderdaad de hashwaarden goed verdeeld zijn, is dan de kans dat er in de emmerketting die herverdeeld moet worden meer dan gemiddeld veel records zitten (de ketting dus langer is) dezelfde als de kans dat er minder dan gemiddeld veel records zitten?*

*Er is geen expliciete berekening vereist, maar geef wel voldoende uitleg waarom jouw antwoord juist is. Bespreek de gevallen dat  $n$  een macht van 2 is – resp. dat niet is – apart.*

Maar belangrijk voor de efficiëntie: maar **één** emmerketting moet herverdeeld worden – en die kan misschien met één enkele leesoperatie ingelezen worden als de hashfunctie en de laadfactor ervoor zorgen dat niet veel overstromemmers nodig zijn. Als je de schrijfoperaties meetelt zijn er dus maar 2 emmerkettingen die gewijzigd moeten worden.

Een voorbeeld van de manier waarop linear hashing werkt, zien jullie in Figuur 8

**Oefening 32** Dit is een bijzonder belangrijke oefening voor linear hashing:

*Stel dat je een goede hashfunctie hebt en dat de reeks van hashwaarden goed verdeeld is en elk van de  $2^b$  mogelijke waarden van  $h()$  ongeveer dezelfde kans maakt om op te duiken. (Om het iets gemakkelijker te hebben mag je stellen dat het aantal sleutels een macht van twee is die groter dan  $2^b$  is en dat elke index **precies** even vaak opduikt.) Zitten de emmers dan altijd allemaal even vol?*

**Als ja:** bewijs dat.

**Als niet:** • Wanneer zijn ze wel even vol?

- Welke emmers zijn voller en hoe groot is het verschil?

*Geef uitleg.*

**Oefening 33** *Geef voldoende voorwaarden op de reeks van sleutels die toegevoegd moeten worden en de hash-functies om te garanderen dat linear hashing voor een reeks van  $n$  toevoegbewerkingen gemiddeld constante tijd vraagt!*  
*Geef uitleg!*

**Oefening 34** Voeg de sleutels 3, 7, 12, 45, 44, 11 en 1 toe aan een lege linear hashing tabel met twee emmers in het begin. Er kunnen 2 sleutels (records) per emmer geplaatst worden en de tabel wordt uitgebreid zodra de laadfactor groter is dan 0.8.

*De hashcodes zijn als volgt:*

$h(3) = 0$  dus binair 0

$h(7) = 31$  dus binair 11111

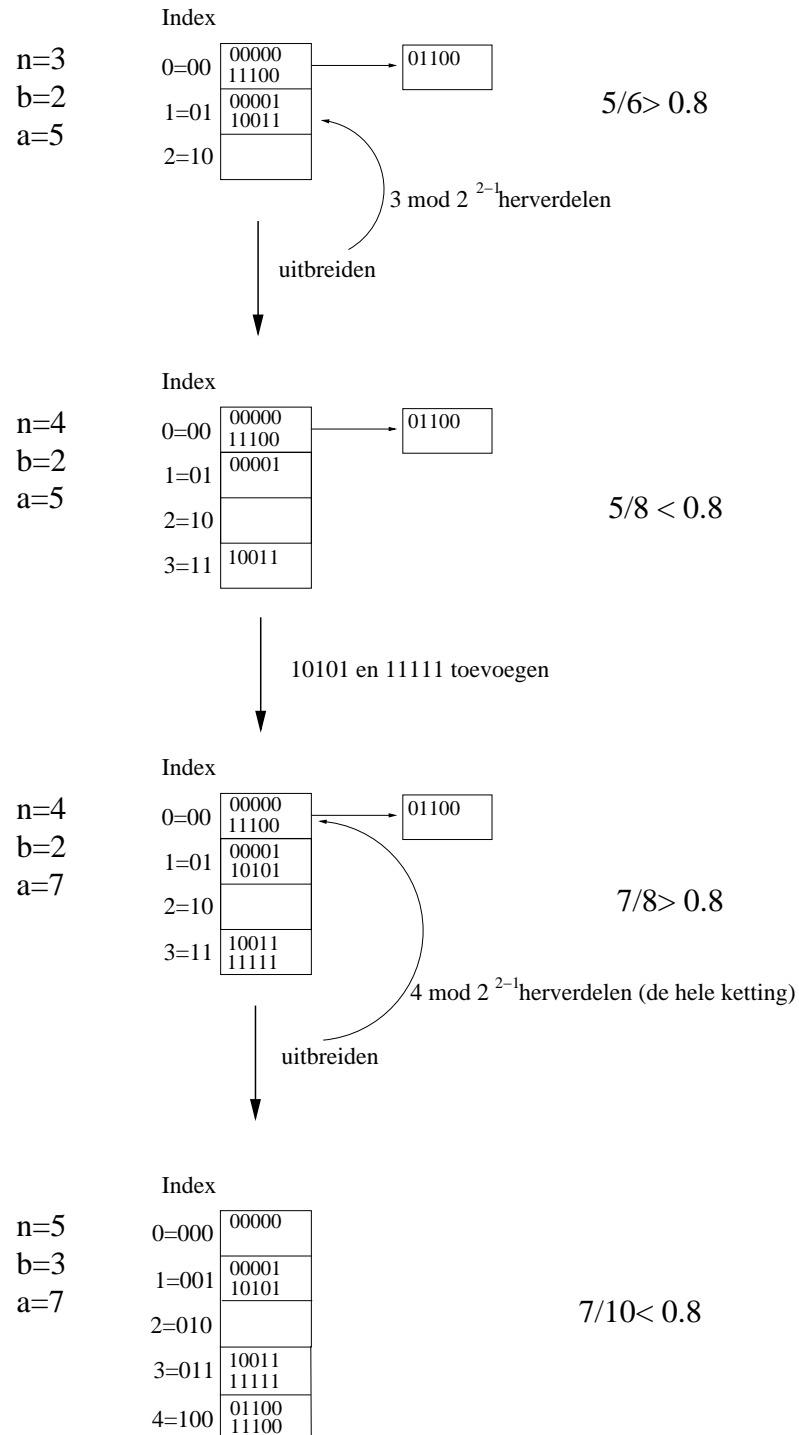
$h(12) = 9$  dus binair 1001

$h(45) = 17$  dus binair 10001

$h(44) = 2$  dus binair 10

$h(11) = 8$  dus binair 1000

$h(1) = 1$  dus binair 1



Figuur 8: Deze linear hashing tabel wordt uitgebreid zodra de laadfactor groter is dan  $L = 0.8$ .

### 3.1.3 Bloom filters

Natuurlijk is het mooi als je een hashing tabel kan uitbreiden als je meer geheugen nodig hebt – maar wat als je gewoon te veel geheugen nodig hebt?

Het volgende algoritme werd in 1970 door Burton Bloom gepubliceerd. Toen was het onmogelijk alle juiste spellingen van woorden in het geheugen op te slaan en het algoritme werd gebruikt om ook met minder geheugen een goede spelling checker te implementeren. Intussen is het beschikbare geheugen zeer veel groter dan toen, maar er zijn nog altijd problemen waar het belangrijk is het geheugengebruik zo klein mogelijk te houden. Soms moet je bv ook informatie over een netwerk sturen en dan is de reductie van de hoeveelheid geheugen die deze informatie codeert ook essentieel! Wij zullen wel zien dat de reductie van geheugengebruik een prijs heeft...

Stel dat je met een universum  $\{0, \dots, n-1\}$  werkt. Als je een grote verzameling van elementen uit dit universum moet bijhouden, dan hebben wij al in Algoritmen en Datastructuren 2 gezien hoe je dat kan doen: je werkt met een bitvector en zet bit nummer  $i$  op 1 als en slechts als element  $i$  in de verzameling zit. Dat vraagt maar 1 bit voor elk element uit het universum. Deze manier van doen werkt natuurlijk ook als de elementen niet  $\{0, \dots, n-1\}$  zijn, maar je wel een unieke index  $\{0, \dots, n-1\}$  aan elk element kan toevoegen.

Eén manier om indices aan objecten toe te kennen is een hashfunctie. Stel dat wij een universum  $U$  van objecten hebben. Dit kunnen bv. alle mogelijke combinaties van ten hoogste 20 lettertekens zijn, maar het kunnen ook alle strings zijn – dus een oneindige verzameling. Stel bovendien dat wij een hashfunctie  $h : U \rightarrow \{0, \dots, n-1\}$  hebben. Als wij nu voor een verzameling  $M \subseteq U$  (bv. alle juist gespelde woorden in de Nederlandse taal) willen toetsen of een zeker element  $x$  in  $M$  zit, kunnen wij als volgt werken:

- Wij nemen een bitvector  $b_0, \dots, b_{n-1}$  met  $n$  bits die in het begin allemaal 0 zijn.
- Dan wordt voor elk element  $y \in M$  de bit  $b_{h(y)}$  op 1 gezet.
- Als wij willen toetsen of een element  $x \in U$  in  $M$  zit, dan berekenen wij  $h(x)$  en zeggen *ja* als bit  $b_{h(x)} = 1$  en anders *nee*.

Als wij dit algoritme toepassen en wij testen een element  $x \in M$  dan zal het antwoord *ja* zijn – dat klopt dus. Als wij een element  $x \notin M$  testen dan kan het antwoord *nee* zijn, maar er is ook een kans dat toevallig een ander element dat wel in  $M$  zit dezelfde hashwaarde heeft en wij dus het (foute) antwoord



*ja* krijgen. Er is dus een kans op false positives, dus foute ja-antwoorden. Je zou de antwoorden als *nee* en *misschien* kunnen interpreteren – dan heb je altijd een juist antwoord, maar aan *misschien* heb je vaak niets... Wij zullen zien hoe je de kans op foute ja-antwoorden kleiner kan maken, maar Bloom filters zijn dus zeker niet geschikt voor gevallen waar een fout ja-antwoord niet achteraf nog door een andere test herkend kan worden (wij zullen een voorbeeld zien waar dat wel kan) en dan misschien rampzalige gevolgen heeft. Als het om een spelling checker gaat dan is een programma dat je een beetje helpt zeker al beter dan helemaal geen hulp – daarvoor waren de Bloom filters dus zeker geschikt. Een andere toepassing waarvoor Bloom-filters voorgesteld werden, was het opslaan van een lijst van onvoldoende veilige passwords. Ook hier was een fout ja-antwoord niet erg: je zou gewoon een ander password moeten kiezen, terwijl foute nee-antwoorden wel erg zouden kunnen zijn – maar foute nee antwoorden heb je bij Bloom filters niet.

Het idee om de kans op false positives te reduceren is eenvoudig: gebruik niet één hashfunctie  $h()$  maar meerdere hashfunctie  $h_1(), \dots, h_k()$  en geef enkel *ja* als antwoord als voor elke hashfunctie  $h_i()$  geldt dat  $b_{h_i(x)} = 1$ .

Daarbij kan je op twee manieren werken als je  $n$  bits ter beschikking hebt:

- a.) Als  $n$  deelbaar is door  $k$  gebruik je voor alle hashfuncties aparte bitvectoren van  $n/k$  bits. De functie  $h_1()$  gebruikt dus bit  $0, \dots, n/k - 1$ , de functie  $h_2()$  gebruikt bit  $n/k, \dots, 2(n/k) - 1$ , etc.
- b.) Je gebruikt voor alle hashfuncties dezelfde  $n$  bits. Het is dus mogelijk dat dezelfde bit door verschillende hashfuncties gezet wordt.

Wij zullen hier alleen maar de (standaard) manier b.) analyseren. Daarbij zullen we de kans berekenen dat een toevallig gekozen element het antwoord *ja* krijgt. Deze kans gebruiken wij hier als benadering voor de kans op een fout positief antwoord. Voor universa die groot zijn in vergelijking met  $M$  is dat zeker een aanvaardbare benadering en voor andere gevallen zou je Bloom filters vermoedelijk toch al niet toepassen...

Op voorhand is daarbij niet duidelijk dat het goed is met meer dan 1 hashfunctie te werken. Als je bv. met extreem veel hashfuncties zou werken, dan zouden bijna alle bits 1 kunnen zijn en je krijgt zeker bijzonder veel foute positieve antwoorden. Is  $k = 1$  misschien zelfs optimaal?

Wij gebruiken dat  $e^x = \lim_{i \rightarrow \infty} (1 + \frac{x}{i})^i$  en dat voor voldoende grote  $i$  geldt dat

$$e^x \approx (1 + \frac{x}{i})^i \tag{1}$$

Bovendien veronderstellen wij dat alle hashfuncties de waarden gelijk verdelen, of precies: wij veronderstellen dat als je een toevallig element  $x \in U$  kiest de kans dat  $h_i(x) = j$  voor alle  $1 \leq j \leq n$  gelijk is aan  $\frac{1}{n}$  en dat voor elke verzameling  $M$  die we beschouwen de kansen voor alle elementen onafhankelijk van elkaar zijn. Dit zijn veronderstellingen die typisch zijn voor hashfuncties en door realistische hashfuncties (en verzamelingen) goed benaderd worden. Wij veronderstellen ook dat de hashfuncties onafhankelijk zijn – dus dat elementen die van de ene hashfunctie een zekere waarde  $i$  krijgen door de andere hashfuncties nog altijd gelijkverdeeld zijn.

Als wij één element met één hashfunctie in onze Bloom filter invullen, dan is de kans dat een zekere bit achteraf nog altijd 0 is dus  $1 - \frac{1}{n}$ . Nadat we  $k$  hashfuncties hebben toegepast, is de kans  $(1 - \frac{1}{n})^k$  en nadat wij  $k$  hashfuncties op  $m = |M|$  elementen hebben toegepast, is de kans  $(1 - \frac{1}{n})^{km}$ .

Als wij in vergelijking 1 voor  $x$  de term  $-\frac{km}{n}$  invullen en voor  $i$  de term  $km$  (waarvan wij hier veronderstellen dat het voor de gekozen  $x$  groot genoeg is), dan krijgen wij als kans dat een bit na het invullen van alle elementen nog altijd 0 is

$$(1 - \frac{1}{n})^{km} \approx e^{-\frac{km}{n}} \quad (2)$$

De kans dat een bit achteraf 1 is, is dus (ongeveer)  $(1 - e^{-\frac{km}{n}})$ .

Omdat wij veronderstellen dat ook de hashfuncties voldoende onafhankelijk zijn, kunnen wij de kansen op een fout positief antwoord schatten: de kans dat een toevallig gekozen element een positief antwoord krijgt (dus dat alle  $k$  bits 1 zijn) is ongeveer  $(1 - e^{-\frac{km}{n}})^k$ .

Als wij de verzameling  $M$  (dus ook  $m = |M|$ ) en het ter beschikking staande geheugen (dus  $n$ ) als gegeven beschouwen, is de vraag hoe we  $k$  moeten kiezen om de kans op foute ja-antwoorden zo klein mogelijk te hebben. Dat is *in principe* heel gemakkelijk: pas jouw kennis uit het middelbaar toe en bereken het minimum van  $f(k) = (1 - e^{-\frac{km}{n}})^k$  waarbij je  $n$  en  $m$  als constanten beschouwt.

Met twee trucjes wordt het iets gemakkelijker:

Wij hebben  $f(k) = (1 - e^{-\frac{km}{n}})^k = e^{k(\ln(1 - e^{-\frac{km}{n}}))}$ . Als wij schrijven

$$g(k) = k(\ln(1 - e^{(-km)/n})) \quad (3)$$

dan heeft  $f(k)$  een minimum als en slechts als  $g(k)$  er één heeft (omdat  $e^x$  strict monotoon stijgt). Als wij nu ook nog schrijven  $p = e^{(-km)/n}$  dan wordt

$$g(k) = g_1(p) = \frac{-n}{m}(\ln p)(\ln(1 - p)). \quad (4)$$

De functie  $g(k)$  heeft dus een minimum als  $g_1(p(k))$  dat heeft en die heeft een minimum als  $g_2(p) = (\ln p)(\ln(1 - p))$  een maximum heeft (omdat wij een negatieve multiplicatieve constante weggelaten hebben).

Daarvoor kunnen wij nu onze kennis uit het middelbaar onderwijs toepassen:

$$\frac{dg_2(p)}{dp} = \frac{\ln(1 - p)}{p} - \frac{\ln p}{1 - p} = \frac{((1 - p) \ln(1 - p)) - (p \ln p)}{p(1 - p)} \quad (5)$$

en dat is 0 voor  $p = \frac{1}{2}$ . Met een beetje meer rekenen kunnen wij ook nog zien dat  $p = \frac{1}{2}$  het enige extremum is en in feite een maximum van  $g_2(p)$ . Met  $p = \frac{1}{2} = e^{(-km)/n}$  krijgen wij dat  $f(k)$  een minimum heeft voor  $k = \frac{n}{m} \ln 2$ .

Als  $m$  en  $n$  gegeven zijn, kiezen we dus het beste  $k = \frac{n}{m} \ln 2$  onafhankelijke hashfuncties om een Bloom filter te implementeren. Natuurlijk moet het aantal een geheel getal zijn, zodat wij dit getal moeten afronden, waarbij wij de voorkeur geven aan  $\lfloor \frac{n}{m} \ln 2 \rfloor$  omdat minder hashfuncties natuurlijk ook minder tijd vragen.

Als je dus bv. een verzameling met 1.000.000 elementen wilt opslaan en je hebt 1 MB ter beschikking, dus 8.000.000 bit, dan gebruik je het best  $\lfloor \frac{8.000.000}{1.000.000} \ln 2 \rfloor = 5$  hashfuncties.

**Oefening 35** • *Stel dat je voor een verzameling  $M$  met  $m$  elementen een Bloom filter wilt maken waar de kans op een fout positief antwoord ten hoogste  $c$  is. Natuurlijk werk je met het optimale aantal hashfuncties. Ontwikkel de formule die het aantal bits beschrijft dat je nodig hebt om een Bloom filter te hebben die aan deze eisen voldoet.*

- *Hoeveel bits heb je nodig om een verzameling met  $m = 1.000.000$  elementen door middel van een Bloom filter voor te stellen als je ten hoogste een kans van 0.02 (dus 2%) op een fout positief antwoord wil hebben?*

Maar het feit dat  $p = \frac{1}{2}$  is ook nog om andere redenen interessant:  $p = e^{(-km)/n}$  is de kans dat een bit bij deze optimale keuze achteraf nog 0 is en omdat  $p = \frac{1}{2}$ , is de kans dat een bit 0 is dus gelijk aan de kans dat de bit 1 is. De structuur van de Bloom filter is dus die van een gelijkverdeelde random string...

**Oefening 36** *Veronderstel nu dat je zoals in a.) werkt en de  $k$  hashfuncties op disjuncte domeinen met grootte  $n/k$  werken. Bereken de kans op foute positieve antwoorden voor gegeven  $k, n, m$  en ook de optimale waarde voor het aantal  $k$  van hashfuncties dat je moet kiezen om de kans op foute positieve antwoorden zo klein mogelijk te maken als  $n, m$  gegeven zijn.*

**Oefening 37** • Als je twee bitmap voorstellingen  $b_M, b_{M'}$  van verzamelingen  $M, M'$  voor hetzelfde universum hebt, dan kan je door de bitsgewijze OR-operator ( $\mid$ ) de verzameling  $M \cup M'$  berekenen:  $M \cup M' = b_M \mid b_{M'}$ . Klopt dat ook voor Bloom filters  $b_M, b_{M'}$  die twee verzamelingen  $M, M'$  voorstellen? Of precies: is het resultaat van  $b_M \mid b_{M'}$  altijd de bloomfilter voor  $M \cup M'$  als het gebruikte aantal bits en de gebruikte hashfuncties dezelfde zijn? Geef uitleg.

- Als je twee bitmap voorstellingen  $b_M, b_{M'}$  van verzamelingen  $M, M'$  voor hetzelfde universum hebt, dan kan je door de bitsgewijze AND-operator ( $\&$ ) de verzameling  $M \cap M'$  berekenen:  $M \cap M' = b_M \& b_{M'}$ . Klopt dat ook voor Bloom filters  $b_M, b_{M'}$  die twee verzamelingen  $M, M'$  voorstellen? Of precies: is het resultaat van  $b_M \& b_{M'}$  altijd de bloomfilter voor  $M \cap M'$  als je hetzelfde aantal bits en dezelfde hashfuncties neemt? Geef uitleg.

**Oefening 38** Stel dat je met een Bloom filter met  $n = 2^i$  bits werkt voor een zekere  $i \in \mathbb{N}, i > 1$ . Je stelt vast dat je  $n$  te groot hebt gekozen en beter met  $n/2$  zou werken, waarvoor het aantal fout positieve antwoorden ook nog aanvaardbaar zou zijn. Beschrijf een manier om een Bloom filter voor  $n/2$  te berekenen zonder alle elementen te rehashen. Welke hash functies gebruik je?

Bloom filters zijn niet alleen nuttig als foute positieve antwoorden niet erg zijn – zoals in het voorbeeld met de slechte passwoorden. Soms is het ook mogelijk achteraf foute positieve antwoorden snel te verbeteren, zoals in het volgende voorbeeld:

In een netwerk heeft een computer  $A$  informatie over alle personen die voor een zeker (groot) bedrijf werken en een andere computer  $B$  informatie over alle adressen in Europa waar er in de toekomst gebouwd moet worden. Computer  $B$  moet nu weten welke personen van het bedrijf op een plaats wonen waar er in de toekomst gebouwd wordt. Eén mogelijkheid zou zijn de hele lijst van plaatsen door te sturen, maar die zou waarschijnlijk te groot zijn. Dus wordt een Bloom filter naar computer  $A$  doorgestuurd waarna computer  $A$  de vermoedelijk kleine lijst van mensen die op zo'n plaats wonen terugstuurt. De weinige *false positives* – dus personen die niet op een plaats wonen waar er gebouwd moet worden – kunnen dan snel verwijderd worden. In dit voorbeeld wordt dus de klemtoon gelegd op weinig data die over het netwerk gestuurd moet worden, omdat dat vaak een bottleneck is.

Natuurlijk is het een voor de hand liggend idee dat je compressiealgoritmen gebruikt om Bloom filters te comprimeren voordat je ze stuurt – ten slotte is weinig geheugen het hoofddoel van Bloom filters. Jammer genoeg hebben

wij net gezien dat in Bloom filters met een minimale kans op een fout positief antwoord elke bit met kans  $p = \frac{1}{2}$  gelijk aan 0 of gelijk aan 1 is – het zijn dus random strings die niet of nauwelijks comprimeerbaar zijn. In 2001 heeft Mitzenbacher iets heel interessants aangetoond: als je jouw Bloom filter niet optimaliseert voor een minimale kans op een fout positief antwoord **voor** het comprimeren, maar **na** het comprimeren, dan kan je soms Bloom filters krijgen die  $n' > n$  bits gebruiken en dezelfde kans op een fout positief antwoord hebben als het optimum voor  $n$  bits – maar beter comprimeren en na compressie minder dan  $n$  bits nodig hebben. Misschien iets om na te lezen. . .

### 3.2 Sorteren van grote hoeveelheden data – extern sorteren

Als je meer data moet sorteren dan je in het geheugen kan houden, kunnen sorteeralgoritmen zoals quicksort die in het geheugen heel goed presteren inefficiënt worden omdat er plotseling een verschil is tussen de verschillende manieren waarop je een element moet lezen – als je een element moet lezen dat al in het geheugen staat is dat veel goedkoper als een element te lezen dat nog van de harde schijf gelezen moet worden. In dit deel zullen wij een algoritme zien dat het aantal leesoperaties op de harde schijf minimaliseert. Onze oplossing zal een mergesort algoritme zijn. Maar in plaats van het in elke recursiestap in maar twee delen te splitsen, zullen wij het op elk niveau in meer delen splitsen waarbij wij rekening zullen houden met de grootte van het geheugen en de grootte van de blokken die je in één keer kan lezen.

**basisstap:** je splitst het hele bestand op zodat je delen  $d_1, \dots, d_n$  hebt die **net** in het geheugen gesorteerd kunnen worden. Deze delen worden efficiënt in het geheugen gesorteerd (bv. door middel van mergesort of quicksort) en teruggestreven.

Als je bv. 10GB moet sorteren en je kan 500MB in het geheugen sorteren (waarvoor je natuurlijk meer dan 500MB geheugen nodig hebt) zonder te moeten swappen, zou je 20 delen maken en die in het geheugen lezen, sorteren en dan terugschrijven. Zo wordt elk blok met data één keer gelezen en één keer geschreven.

**samenvoegstap:** Nu moeten de gesorteerde delen nog gemerged worden. Wij lezen van zoveel mogelijk delen het eerste blok in het geheugen. *Zoveel mogelijk* betekent dat wij maar zoveel delen kiezen dat wij voor elk deel één blok in het geheugen kunnen houden. Dan schrijven wij altijd de kleinste sleutel van één van de delen in een outputblok tot dat

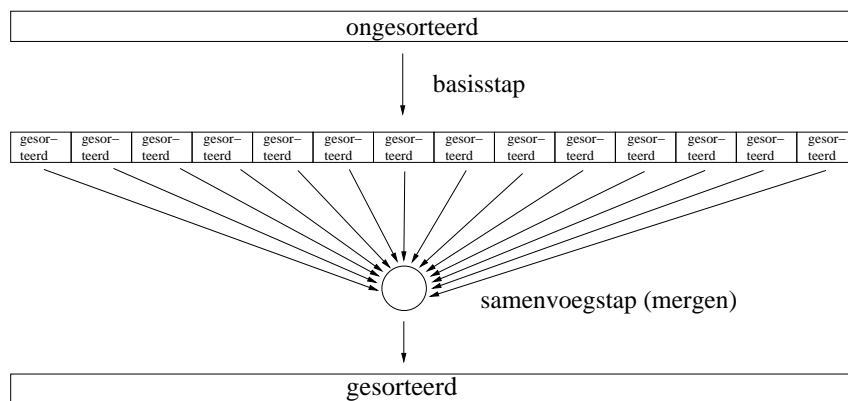
volzit en uitgeschreven wordt. Als één van de blokken leeg is, wordt dat door het volgende blok van dat deel vervangen.

Als wij  $t$  delen tegelijk kunnen mergen en er zijn samen  $d$  delen dan passen wij dit  $\lceil d/t \rceil$  keer toe totdat elk deel één keer met andere delen samengemerged werd. Wij hebben dan  $\lceil d/t \rceil$  grotere gesorteerde delen.

In deze stap wordt elk blok opnieuw één keer gelezen en één keer geschreven.

Deze samenvoegstap wordt herhaald tot er maar één bestand overblijft.

Het principe voor één samenvoegstap zien jullie nog eens in Figuur 9 en in Figuur 10 zien jullie een voorbeeld met een heel kleine vertakking van 4. Daardoor zijn er twee samenvoegstappen nodig. Realistische voorbeelden kan je natuurlijk niet tekenen. . .



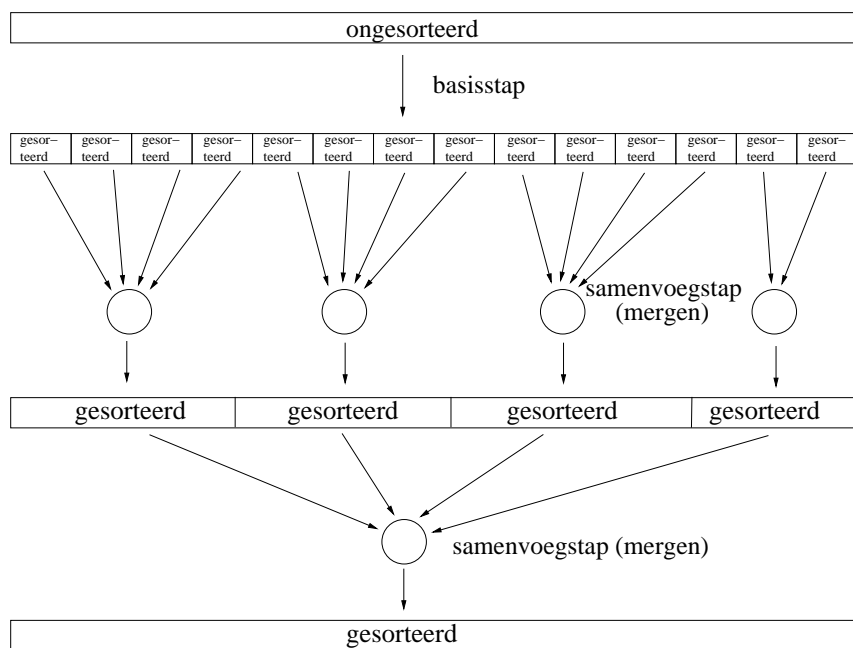
Figuur 9: Het principe van extern sorteren.

Stel dat  $B$  de blok grootte is die je in één stap kan lezen en  $M$  de grootte van het geheugen dat je voor het opslaan van de delen ter beschikking kan stellen. Dan kan je dus  $\frac{M}{B}$  delen per samenvoegstap mergen. Omdat in het begin de delen grootte  $M$  mogen hebben, kan je dus in  $s$  samenvoegstappen bestanden van grootte  $M * (\frac{M}{B})^s$  sorteren.

Als je er opnieuw van uitgaat dat je 500MB ter beschikking hebt en dat bv. een blok van 4KB in één stap gelezen kan worden, zijn dat al 62.5 Terabyte die je kan sorteren met maar één enkele samenvoegstap. Dat is inderdaad heel veel en ook al wordt de hoeveelheid data groter – het geheugen groeit ook en de 500 MB waarvan wij hier gesteld hebben dat die ter beschikking staan, zijn zeker een benedengrens. . .

Wij kunnen dus met 2 lees- en schrijf-operaties per blok 62.5 TB sorteren.

Wat wij hier gezien hebben is dus gewoon mergesort waar wij aan de basis niet alleen enkele sleutels hebben maar al grotere reeds gesorteerde delen.



Figuur 10: Extern sorteren in 2 stappen. In dit voorbeeld kunnen maar 4 delen tegelijk gemerged worden.

En het belangrijkste: de vertakking is **veel** groter. De vertakking  $\frac{M}{B}$  bepaalt de diepte van de recursieboom en de diepte van de recursieboom bepaalt het aantal lees- en schrijf-operaties.

Wij hebben er niet op gelet wat je bv. moet doen als de records die je wil sorteren niet in één blok passen etc. In dergelijke gevallen zijn lichte wijzigingen nodig.

**Oefening 39** *Wij hebben geschreven „Dan schrijven wij altijd de kleinste sleutel van één van de delen in een outputblok...“. Hoe doe je dat het best op een efficiënte manier? Altijd naar het kleinste element in elk opgeslagen blok kijken zou lineair zijn in het aantal blokken. Kan dat efficiënter? Denk aan DA1 en DA2.*

**Oefening 40** *Gegeven de grootte  $M$  van het geheugen dat ter beschikking staat en  $B$  de blok grootte die je in één stap kan lezen. Geef een formule voor het aantal leesoperaties (per blok) op de harde schijf dat nodig is om een bestand met  $n$  bytes te sorteren.*

**Oefening 41** *De bedoeling van deze oefening is (nog) een beetje meer gevoel te krijgen voor welke betekenis de asymptotische analyse heeft en in welke gevallen ze meer of minder belangrijk is.*

*Stel dat een lees/schrijf-operatie op de harde schijf kost  $C_1$  per blok heeft en een operatie in het geheugen (vergelijken, verplaatsen, etc.) kost  $C_2$ . De motivatie voor de laatste algoritmen die wij hebben gezien was dat  $C_1 \gg C_2$  – dus mag je dat ook hier stellen.*

*Je hebt een vast geheugen van 400 MB ter beschikking en een bestand van  $n$  records waarvan elke record een grootte van 400 bytes heeft. Een blok heeft grootte 4 KB.*

*Onderzoek drie varianten van extern sorteren: één keer sorteer je de deelbestanden die net in het geheugen passen met mergesort, één keer met bubble-sort en één keer met quicksort. In elk geval pas je achteraf de mergebewerking op zoveel niveau's toe als er noodzakelijk zijn.*

- *Zal er een groot verschil in performantie zijn als je deze 3 varianten in de praktijk op een bestand van bv. 10 GB (dus 25000000 records) toepast?*
- *Bepaal de asymptotische complexiteit van alle drie varianten van extern sorteren.*

### **3.3 Grote zoekbomen**

Waarschuwing in het begin: in de literatuur vind je B-trees en B+-trees soms op verschillende manieren gedefinieerd. De ideeën zijn natuurlijk altijd dezelfde als jullie ze hier zullen zien, maar de details zijn soms een beetje verschillend – dus opgelet als jullie op verschillende plaatsen kijken.

In dit deel zullen wij het nu over een echte vertakking hebben – het gaat over bomen. In zoekbomen is het aantal kinderen van een top altijd één groter dan het aantal sleutels die de top bevat (waarbij de kinderen leeg kunnen zijn). Als wij een boom hebben waar het kleinste aantal niet lege kinderen van een top die geen blad is  $g$  is (elke interne top moet dan ten minste  $g - 1$  sleutels bevatten en stel dat dat ook voor de bladeren geldt) en alle bladeren zitten op dezelfde diepte  $d$  dan kunnen in deze boom ten minste  $g^{d+1} - 1$  sleutels geplaatst worden. Om  $n$  sleutels te plaatsen, heb je dus een diepte van ongeveer  $\log_g(n) = \frac{1}{\log g} \log n$  nodig. Als de diepte het aantal leesoperaties beschrijft, moeten wij er dus voor zorgen de graad van de toppen zo groot mogelijk te kiezen. En precies dat zullen wij doen: wij proberen een boom met een vertakking te bouwen die zo groot mogelijk is – onder de voorwaarde dat wij elke top nog met een enkele leesoperatie kunnen lezen en de top dus nog in één blok geplaatst kan worden.



### 3.3.1 B-trees

Wij herhalen in de definitie van een B-tree de ordeningseigenschappen in de toppen van de boom in plaats van gewoon te zeggen dat het een zoekboom is, zodat de definitie beter met die van een B+-tree vergeleken kan worden.

**Definitie 1** *Wij stellen hier dat een boom een sleutel maar één keer kan bevatten.*

*Een B-tree met (even) grootte  $n$  is een boom met volgende eigenschappen:*

- *het aantal  $a$  van sleutels in de wortel voldoet aan  $1 \leq a \leq n$*
- *het aantal  $a$  van sleutels in toppen die niet de wortel zijn, voldoet aan  $\frac{n}{2} \leq a \leq n$*
- *als in een top  $t$  de sleutels  $s_1 < s_2 < \dots < s_a$  zitten dan geldt*
  - *Ofwel is  $t$  een blad (en heeft dus geen kinderen) ofwel heeft  $t$  precies  $a + 1$  (niet lege) kinderen  $k_1, \dots, k_{a+1}$  waarvoor met  $s_0 = -\infty$  en  $s_{a+1} = \infty$  voor alle sleutels  $s$  in de deelboom met wortel  $k_i$  geldt dat  $s_{i-1} < s < s_i$ .*
- *alle bladeren zitten op dezelfde afstand van de wortel (dezelfde diepte).*

Het is dus duidelijk een uitbreiding van de definitie van een 2-3-boom en ook de bewerkingen lijken heel sterk op die van een 2-3-boom:

**zoeken:** Voor elke top zoek je of de sleutel in die top zit of in welk kind je moet doorgaan met het zoeken. Voor grote  $n$  doe je dat het best door middel van logaritmisch zoeken.

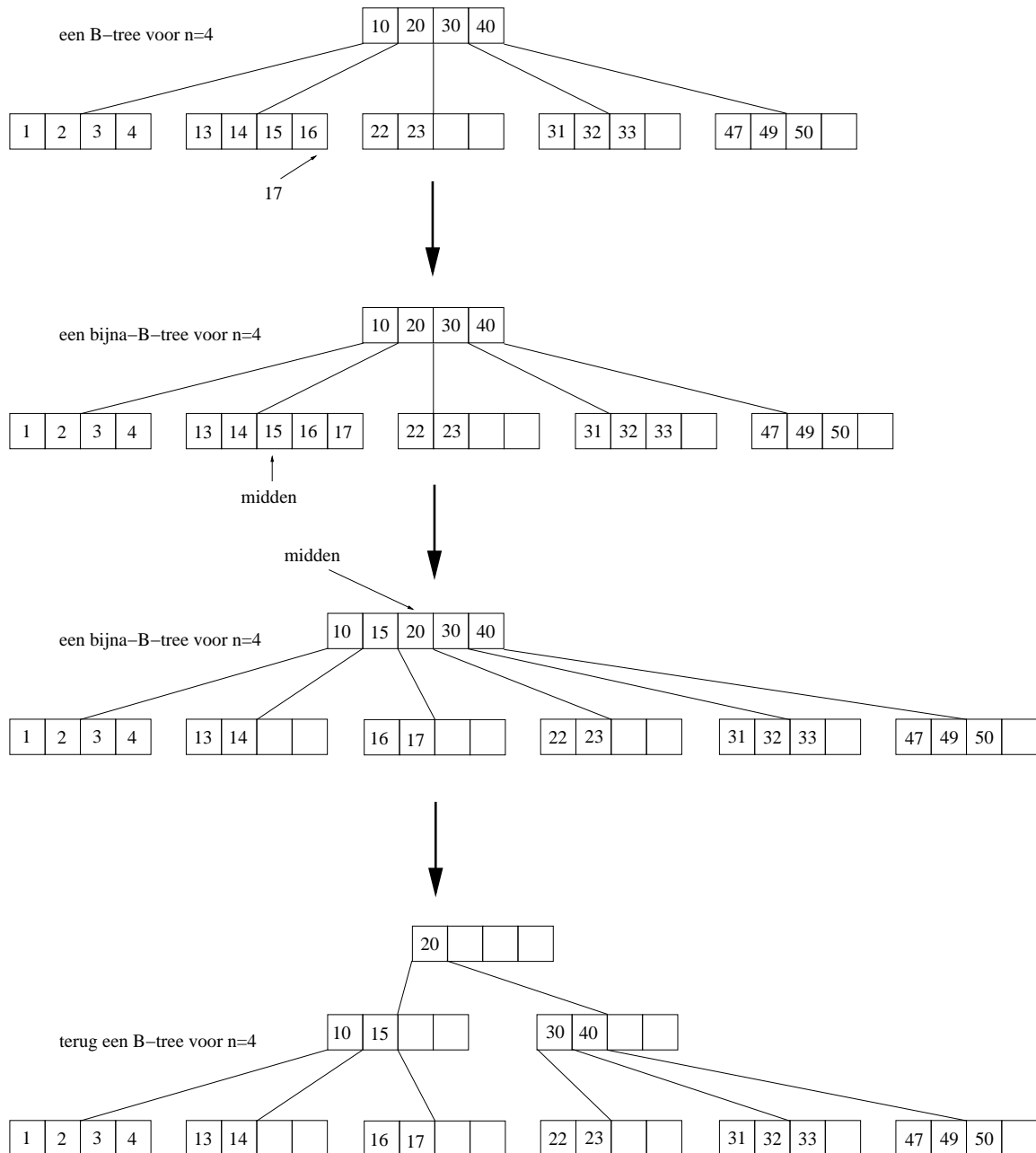
**toevoegen:** Eerst zoek je het blad waar de sleutel toegevoegd moet worden.

- Als er nog plaats is, voeg je hem gewoon toe (waarbij de sleutels in het blad gesorteerd moeten blijven).
- Anders bevat de top al  $n$  sleutels, dus samen met de nieuwe  $n + 1$  sleutels. Wij werken nu tijdelijk met een boom waar **één** top  $n + 1$  sleutels bevat (dat noemen wij hier een bijna-B-tree) en die herbalanceren wij dan tot hij opnieuw aan de eigenschappen van een B-tree voldoet.

Het principe is daarbij gemakkelijk om te verstaan: omdat  $n + 1$  sleutels niet in één top passen, heb je er twee nodig – de ouder zal dus één kind meer hebben. Maar omdat het aantal kinderen van de ouder bepaald wordt door

het aantal sleutels, moet ook één sleutel naar boven schuiven om ervoor te zorgen dat er één kind meer is.

Een voorbeeld voor het herbalanceren zien jullie in Figuur 11.



Figuur 11: Het toevoegen en herbalanceren in een B-tree voor  $n = 4$ .

Daarbij moet je de bijna-B-tree vooral als een model zien. Jullie moeten in de implementatie niet echt met een top met  $n + 1$  sleutels werken – jullie kunnen ook gewoon met een B-tree werken waar je voor één top nog een extra sleutel hebt. Het model waar er  $n + 1$  sleutels in één top zitten, is gewoon gemakkelijker om uit te leggen. Het herbalanceren gebeurt precies als volgt: Je splitst de top met  $n + 1$  sleutels. De grootste  $n/2$  sleutels zitten in de ene nieuwe top en de kleinste  $n/2$  sleutels in de andere. Dan heb je nog de middelste sleutel  $s$  die ertussen zit. Als er geen ouder is, wordt die de nieuwe wortel met deze twee toppen als kinderen (voor de wortel is het toegelaten dat hij maar 1 sleutel bevat – om het even wat  $n$  is). Als er een ouder is wordt  $s$  als sleutel aan de ouder toegevoegd en de twee nieuwe toppen als kinderen. Omdat de ouder nu één sleutel meer bevat moet hij ook één kind meer hebben. De positie waar je de sleutel en de kinderen moet toevoegen vind je snel als je nog weet het hoeveelste kind de top was. Op deze manier hebben wij één top die te veel sleutels bevatte verwijderd – maar misschien ook één nieuwe top gemaakt (als de ouder al  $n$  sleutels bevatte). Maar die zit nu dicht bij de wortel – als wij het herstellen recursief toepassen, gaat het proces dus zeker stoppen. Wij moeten alleen maar naar toppen langs het toegangspad kijken – dat is dus goedkoop – en bovendien hebben wij nadat een top gesplitst moet worden twee toppen met samen  $n$  vrije plaatsen – dus voldoende ruimte voor  $n$  toevoegoperaties zonder splitsen. Samenvattend kan dus gezegd worden dat toevoegen een bewerking is die heel efficiënt is.

### **Verwijderen:**

Net zoals voor 2-3-bomen is verwijderen ook hier relatief ingewikkeld en duur. Als er niet te veel verwijderbewerkingen zijn, is het dus het beste als je met grafstenen werkt: als een sleutel verwijderd moet worden, wordt hij niet echt verwijderd maar gemarkeerd als niet meer bestaande. Je zou hem natuurlijk ook kunnen vervangen door de grootste sleutel in zijn kleiner-kind of de kleinste sleutel in zijn groter-kind en inderdaad alleen sleutels in de bladeren verwijderen. Alleen als die beide ook een grafsteen hebben zou je ook grafstenen in het midden van de boom plaatsen. In een interne top moet de sleutel aanwezig blijven omdat hij tijdens het zoeken en toevoegen ook nog gebruikt wordt voor de *navigatie* in de boom – dus om de juiste sleutels en plaatsen te vinden. In een blad zou je hem ook gewoon kunnen verwijderen, zodat *grafstenen* in principe betekent dat je hier toppen met minder dan  $n/2$  sleutels aanvaardt.

**Oefening 42** *Werk precies uit hoe je in een B-tree met grafstenen werkt, dus wanneer je echt grafstenen plaatst of gewoon toelaat dat er minder sleutels zijn, wanneer je een sleutel met een grafsteen door een nieuwe sleutel*

vervangt, etc.

Als het een boom is met **extreem** veel verwijderbewerkingen (zonder tussendoor voldoende toevoegbewerkingen die de grafstenen in de bladeren toch automatisch verwijderen) zou je misschien beter toch niet altijd met grafstenen werken. Wij zullen er ons hier gewoon van overtuigen dat je op een manier die alleen lokale wijzigingen langs een pad van de wortel naar een blad gebruikt de boom **kan** herbalanceren zonder de details uit te werken (maar dat kan je dan wel zelfstandig). Maar dat is vooral omdat wij nieuwsgierig zijn of het kan...

**Oefening 43** *Bewijs het volgende lemma:*

**Lemma 2** *Gegeven een even getal  $n \geq 2$ . Dan bestaat er een B-tree met grootte  $n$ , ten minste  $n/2$  sleutels in de wortel, met  $s$  sleutels en diepte 1 als en slechts als  $\frac{n^2}{4} + n \leq s \leq n^2 + 2n$ .*

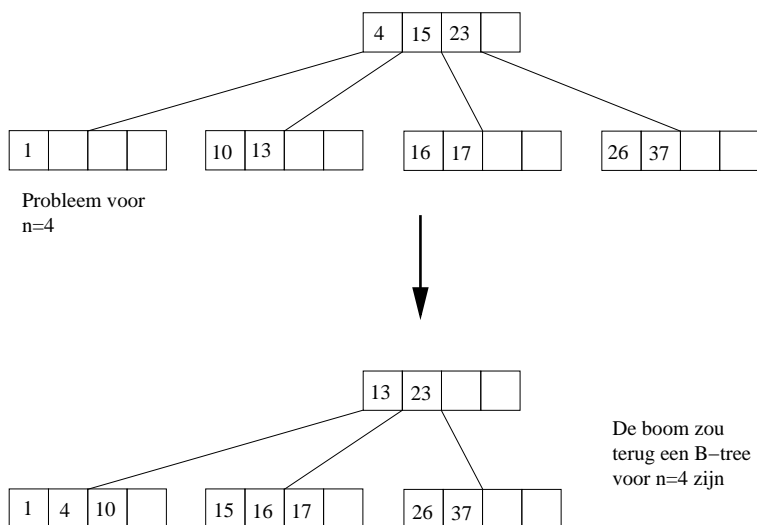
Zoals net gezegd kunnen wij er ook hier vanuitgaan dat wij alleen maar sleutels in bladeren verwijderen. Als een sleutel uit een top wordt verwijderd die geen blad is, kunnen wij die vervangen door de kleinste sleutel in de deelboom waarvan het *groterkind* de wortel is of de grootste sleutel in de deelboom waarvan het *kleinerkind* de wortel is. Die zitten zeker in een blad. Als het blad achteraf nog ten minste  $\frac{n}{2}$  sleutels bevat, moeten wij niets doen. Anders werken wij met een bijna-B-tree waar deze keer één top één sleutel te weinig bevat.

Stel eerst dat de ouder van de foutieve top niet de wortel is. Dan zitten er in de ouder en de kinderen van de ouder samen ten minste  $\frac{n^2}{4} + n - 1$  sleutels en ten hoogste  $n^2 + \frac{3}{2}n - 1$  sleutels. Dus kan – behalve in het geval dat er precies  $\frac{n^2}{4} + n - 1$  toppen zijn – de deelboom bestaande uit de ouder en zijn kinderen vervangen worden door een andere deelboom die aan de eisen voldoet (deelboomvervangmethode en Lemma 2). Achteraf is de boom opnieuw een B-tree.

Als er precies  $\frac{n^2}{4} + n - 1$  sleutels inzitten, gaan wij de deelboom vervangen door een deelboom waar de ouder maar  $\frac{n}{2} - 1$  sleutels heeft. De foutieve top zit dan één stap dicht bij de wortel en door de hersteloperaties herhaaldelijk toe te passen, kunnen wij de boom uiteindelijk weer tot een B-tree maken. Om te bewijzen dat je inderdaad deze iets te kleine deelbomen ook door deelbomen met de fout in de wortel kan vervangen, kunnen wij (iets algemener dan echt nodig) analoog met Lemma 2 aantonen dat er voor  $\frac{n^2}{4} + \frac{n}{2} - 1 \leq s \leq \frac{n^2}{2} + \frac{n}{2} - 1$  sleutels zo'n bomen met  $\frac{n}{2} - 1$  sleutels in de wortel bestaan en wij hebben duidelijk  $\frac{n^2}{4} + n - 1 \geq \frac{n^2}{4} + \frac{n}{2} - 1$  en  $\frac{n^2}{4} + n - 1 \leq \frac{n^2}{2} + \frac{n}{2} - 1$  kan voor  $n \geq 1$  ook gemakkelijk bewezen worden. De nodige bomen bestaan dus.

Wat overblijft is het geval dat de ouder de wortel is. Dit is volledig analoog – behalve dat wij in dit geval aan de ene kant ook naar het geval moeten kijken dat er minder dan  $\frac{n}{2}$  sleutels in de ouder zitten (in de wortel mag dat) maar er aan de andere kant geen rekening mee moeten houden dat er achteraf ten minste  $\frac{n}{2}$  sleutels in de ouder moeten zijn. Het enige geval waar wij de deelboom hier niet door een deelboom met diepte 1 kunnen vervangen, is waar er maar  $n$  sleutels zijn. In dit geval nemen wij een boom met maar 1 top en  $n$  sleutels.

Natuurlijk zijn de details hier niet echt uitgewerkt maar ik hoop toch dat jullie ervan overtuigd zijn dat de boom hersteld kan worden door alleen maar lokale bewerkingen langs het toegangspad te doen. Jammer genoeg kan ook een lokale bewerking hier veel meer leesoperaties vragen dan wij willen investeren – dit is dus alleen een oplossing als grafstenen een veel te grote overhead zouden betekenen. Als je de verwijderoperatie **echt** wil implementeren moet je natuurlijk goed over de meest efficiënte manier nadenken en dan zal je zeker niet onmiddellijk naar alle kinderen van een top kijken maar eerst of je het probleem niet al met een deel van de kinderen kan oplossen. . .



Figuur 12: Een voorbeeld voor het vervangen van een deelboom waar één top één sleutel te weinig bevat. De bijna-B-tree waarin dit deel wordt vervangen, zou achteraf terug een B-tree zijn.

Wij hebben B-trees geïntroduceerd om ervoor te zorgen dat de diepte – die in ons model het aantal leesoperaties beschrijft – zo klein mogelijk is. Maar hoe groot is de diepte ten hoogste als je  $s$  sleutels hebt?

Het slechtste geval (dus het geval waar voor een gegeven diepte het aantal sleutels minimaal is) is duidelijk als de wortel maar één sleutel bevat en

alle andere toppen maar  $\frac{n}{2}$  sleutels. Als de diepte  $d \geq 1$  is en schrijven wij  $t = \frac{n}{2} + 1$  dan bevat elk van de twee kinderen van de wortel

$$(t - 1)(1 + t + t^2 + t^3 + \dots + t^{d-1}) = t^d - 1$$

sleutels. Samen met de wortel zijn dat  $2t^d - 1$  sleutels. In een B-tree van orde  $n$  met  $s$  sleutels is de diepte  $d \leq \log_t \frac{s+1}{2} = (\log \frac{s+1}{2}) / \log(\frac{n}{2} + 1)$ .

### Hoe kies je $n$ ?

Wij willen  $n$  zo groot mogelijk kiezen, maar het moet wel nog altijd in één leesbewerking gelezen kunnen worden.

Stel bv. dat wij getallen met 5 bytes willen opslaan en een blok heeft 4096 bytes. Wij hebben ook nog de pointers nodig. Stel dat die 8 bytes nodig hebben. Voor  $n$  sleutels per top hebben wij dus  $5n + 8(n + 1)$  bytes nodig. Dus kiezen wij een even  $n$  zo groot mogelijk maar op een manier dat nog steeds  $5n + 8(n + 1) \leq 4096$ . Dat geeft  $n = 314$ .

Als wij in een B-tree met deze parameters  $10^{10}$  getallen moeten plaatsen, hebben wij dus een B-tree met diepte  $d \leq (\log \frac{10^{10}+1}{2}) / \log(\frac{314}{2} + 1) < 4.411$ . Dus  $d = 4$  in het slechtste geval. Dus zijn ten hoogste 5 leesoperaties nodig om elke van de  $10^{10}$  sleutels te vinden!

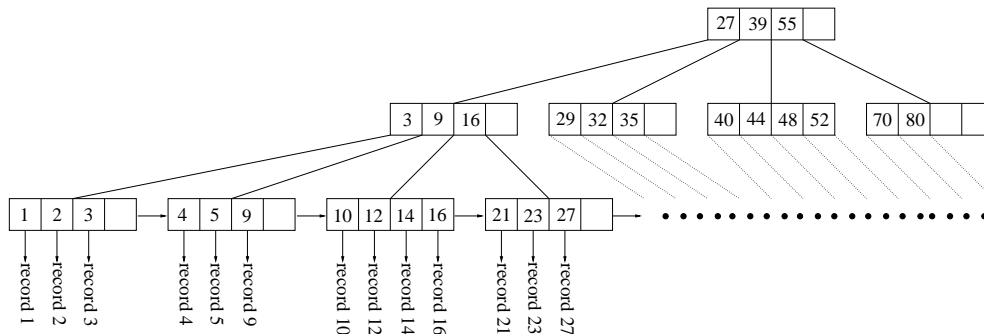
**Oefening 44** *Stel dat 99% van de sleutels in een B-tree met grootte 300 grafstenen hebben en  $s$  het aantal sleutels zonder grafstenen is. Hoeveel slechter is de bovengrens voor de diepte van deze B-tree in vergelijking met een B-tree met grootte 300 waarin alleen maar de  $s$  sleutels zonder grafstenen zitten?*

**Oefening 45** *Stel dat je een B-tree met grootte  $n$  hebt waar in elke top precies  $n$  sleutels zitten.*

- *Hoe groot is het aantal sleutels in de bladeren in vergelijking met het aantal sleutels in interne toppen van de boom?*
- *Stel dat er  $s$  sleutels in de bladeren zitten (onder deze omstandigheden moet dus  $s = n * (n + 1)^d$  voor een  $d$  zijn). Bewijs een goede bovengrens voor het verschil in diepte tussen deze boom en een boom  $B'$  waar alleen maar de sleutels uit de bladeren inzitten. Let op: voor  $B'$  kan je natuurlijk niet eisen dat er precies  $n$  sleutels in elke top zitten – kijk één keer naar de minimale en één keer naar de maximale diepte.*

### 3.3.2 B+-trees

Maar hoe realistisch is het dat wij met getallen als *typische* gegevens werken? Normaal is het zeker zo dat de records met de hele data veel groter zijn. De data ter identificatie (naam, rijksregisternummer, etc.) is misschien wel relatief klein, maar bovendien kan de record nog veel meer informatie bevatten (misschien zelfs een foto of andere gegevens die veel geheugen vragen zoals contracten, facturen, etc.). Deze informatie in de records zou de mogelijke vertakking **extreem** beperken. Het zou dus een goed (en voor de hand liggend) idee zijn niet echt de records op te slaan maar gewoon sleutels die nodig zijn om de records te identificeren en een pointer naar een record op de harde schijf – de boom dus gewoon als een index te gebruiken. Dat zou je dus zeker ook in *B*-trees doen – anders wordt de vertakking veel te klein. Maar B+-trees gaan nog één stap verder: ook de extra pointer vraagt natuurlijk geheugen en beperkt de vertakking. Daarom zullen wij pointers naar de records **alleen** in de bladeren bijhouden. Daar heb je de ruimte voor de pointers naar de kinderen niet nodig (als je in 1 bit of 1 byte bijhoudt dat het een blad is of – nog beter – gewoon bijhoudt wat de diepte van de boom is) en kan je dezelfde ruimte misschien gebruiken voor de pointer naar de records. Dat betekent wel dat elke sleutel in de bladeren aanwezig moet zijn, dus dat er kopieën in de binnenste toppen zitten – en wij dus meer sleutels hebben. Maar door de grote vertakking is de verhouding tussen het aantal bladeren en het aantal interne toppen heel groot.



Figuur 13: Het principe van een B+-tree.

In Figuur 13 zien jullie het principe van een B+-tree. De pointers die van het ene blad naar het andere blad gaan, maken de boom natuurlijk tot iets dat inderdaad geen boom is – de naam is dus een beetje misleidend. De bedoeling is efficiënter van het ene blad naar het volgende te kunnen gaan – bv. als je over alle gegevens wil itereren.

**Definitie 2** Wij stellen hier dat elke sleutel in maar één record voorkomt. Een B+-tree met (even) grootte  $n \geq 2$  is een boom met volgende eigenschappen:

- het aantal  $a$  van sleutels in de wortel voldoet aan  $1 \leq a \leq n$
- het aantal  $a$  van sleutels in toppen die niet de wortel zijn, voldoet aan  $\frac{n}{2} \leq a \leq n$
- als in een top  $t$  sleutels  $s_1 < s_2 < \dots < s_a$  zitten dan geldt
  - Ofwel is  $t$  een blad (en heeft dus geen kinderen) ofwel heeft  $t$  precies  $a + 1$  kinderen  $k_1, \dots, k_{a+1}$  waarvoor met  $s_0 = -\infty$  en  $s_{a+1} = \infty$  voor alle sleutels  $s$  in de deelboom met wortel  $k_i$  geldt: **(let op – gewijzigd in vergelijking met B-trees:)** dat  $s_{i-1} < s \leq s_i$ .
- alle bladeren zitten op dezelfde afstand van de wortel (dezelfde diepte).
- **nieuw:** Elke sleutel van een record zit precies één keer in een blad.
- **nieuw:** Als  $t$  een blad met  $a$  sleutels  $s_1, \dots, s_a$  is, bevat  $t$  ook  $a$  verwijzingen  $p_1, \dots, p_a$  naar records waarbij voor  $1 \leq i \leq a$  geldt dat  $p_i$  naar de record wijst die sleutel  $s_i$  bevat.

Bovendien is er nog een extra pointer die de eigenschap een boom te zijn verstoort (maar wij noemen het toch al B+-tree):

- **nieuw:** Elk blad bevat een pointer naar het volgende blad. (De pointer van het laatste blad wijst naar NULL.)

Het maakt in principe niet uit of je in plaats van  $s_{i-1} < s \leq s_i$  vraagt dat  $s_{i-1} \leq s < s_i$ . Het resultaat zal een even efficiënte datastructuur zijn en jullie zullen beide definities in de literatuur terugvinden.

Je zal ook definities vinden waar de pointers tussen de bladeren in beide richtingen gaan – dus zodat het volgende en het vorige blad bereikt kunnen worden.

De bewerkingen lijken heel sterk op die van B-trees. Maar er zijn natuurlijk sommige verschillen die veroorzaakt zijn door het feit dat elke sleutel in een blad moet opduiken.

**zoeken:** Zoeken is inderdaad het gemakkelijkst. Stel dat je sleutel  $s$  zoekt en met  $s_i$  vergelijkt. Je moet hier dan zelfs niet testen of  $s = s_i$  als het een interne top is omdat je in de gevallen *gelijk* en *kleiner* op dezelfde



manier moet doorgaan – je moet in één van de kinderen zoeken die een index van ten hoogste  $i$  hebben. Het juiste kind bepalen kan je ook hier bv. met logaritmisch zoeken.

**toevoegen:** Als je een sleutel in een blad kan toevoegen zonder dat het volzit, is er geen probleem. Als dat wel volzit moet je splitsen: je splitst de top in twee toppen met een grootte die zo gelijk mogelijk is (dan verschilt het aantal sleutels in de twee toppen ten hoogste met 1). Let op de pointers tussen de bladeren! Nu heeft de ouder één kind te veel en moet hij één sleutel meer krijgen. Je neemt een kopie van de grootste sleutel in de kleinere van de twee delen en voegt die aan de ouder toe. Nu kan je er ook een pointer naar een kind meer plaatsen. Als de ouder niet volzat is het gedaan – anders heb je nu één top met één sleutel te veel. Dan wordt op dezelfde manier doorgegaan als in het geval van B-trees. Zie Figuur 14 voor een voorbeeld.

**verwijderen:** Ingewikkeld – maar jullie hebben nu al voldoende kennis om het verwijderen analoog met B-trees zelf uit te werken. Maar ook hier geldt dat als er niet extreem veel verwijderbewerkingen zijn het gewoon beter is grafstenen te gebruiken.

**Oefening 46** *In een top mogen geen twee identieke sleutels zitten. Maar wij kopiëren een sleutel uit een blad en voegen die aan de ouder toe – hoewel het wel kan gebeuren dat een sleutel in een blad en zijn ouder zit. Toon aan dat de gekopieerde sleutel nog niet in de ouder van het blad zit.*

B+-trees kunnen dus even gemakkelijk gebruikt worden als B-trees. Toevoegen en opzoeken kan heel efficiënt geïmplementeerd worden en alleen de verwijderoperatie is een beetje ingewikkeld. Het feit dat de interne toppen in principe alleen maar voor de navigatie in de boom dienen en je kopieën van sleutels hebt, wordt door het feit dat je geen pointers naar records in de interne toppen moet bijhouden meer dan gecompenseerd.

Maar hoe groot is de diepte van een B+-tree met grootte  $n$  in het slechtste geval? Het geval waar voor een gegeven diepte het aantal sleutels minimaal is, is ook hier duidelijk als er 1 sleutel in de wortel en  $\frac{n}{2}$  sleutels in alle andere toppen zitten. Als de diepte  $d \geq 1$  is en schrijven wij  $t = \frac{n}{2} + 1$  dan bevat elk van de twee kinderen van de wortel  $t^{d-1}$  bladeren dus verwijzingen naar  $(t-1)t^{d-1}$  records. In de hele boom kunnen dus  $2(t-1)t^{d-1}$  records opgeslaan worden, dus  $s \geq 2(t-1)t^{d-1}$ . Als wij dat herschrijven dan geldt voor de diepte  $d \leq \log_t s + 1 - \log_t(2(t-1))$  en omdat  $2(t-1) \geq t$  geldt  $\log_t(2(t-1)) \geq 1$  (maar het is nauwelijks groter) en  $d \leq \log_t s = \frac{\log s}{\log(\frac{n}{2}+1)}$ .

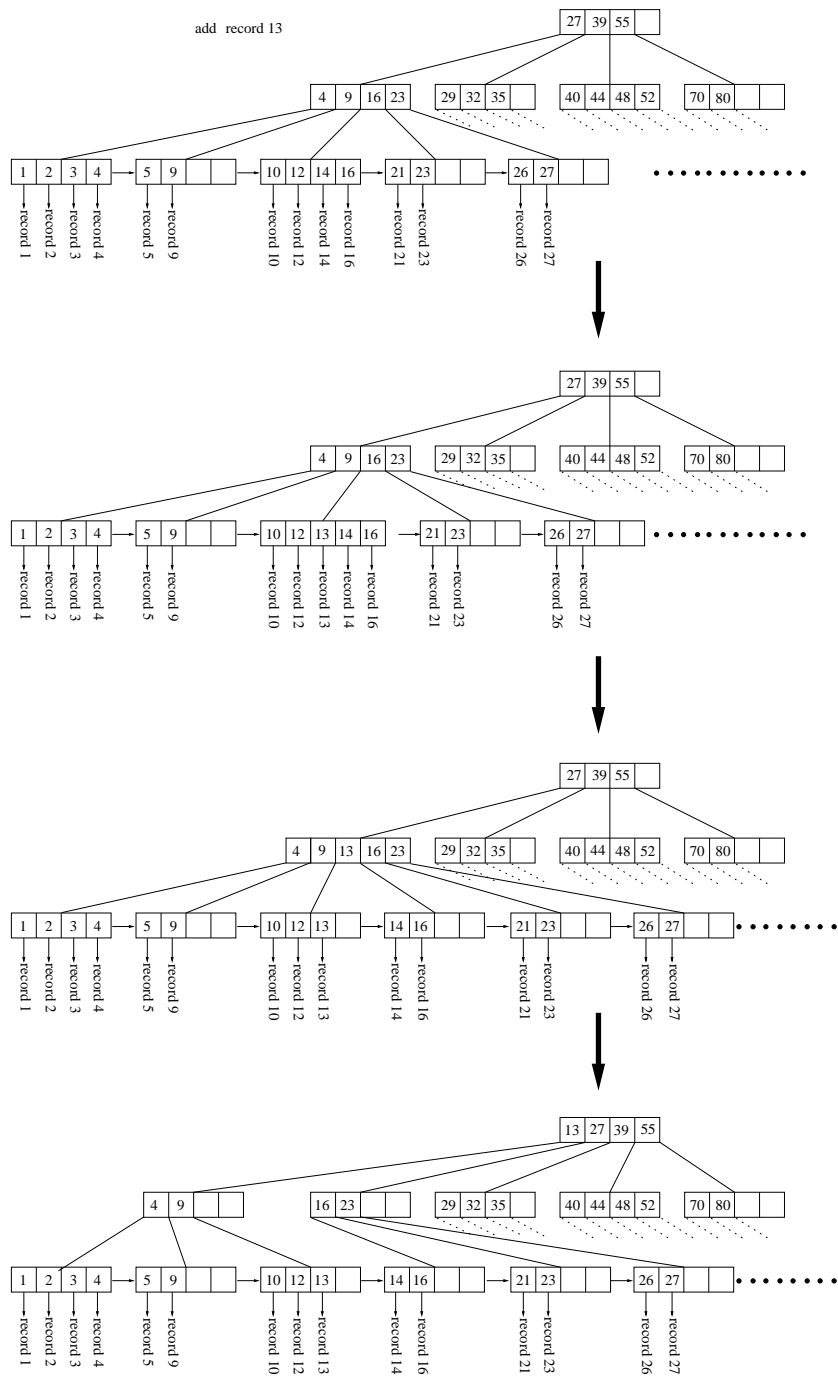
Als wij dat nu met de grens voor B-trees vergelijken dan zien wij dat dat zelfs voor dezelfde vertakking  $t$  naar beneden afgerond meestal dezelfde diepte is en alleen in zeldzame gevallen 1 groter. En als je er nog rekening mee houdt dat de vertakking voor B-trees groter gekozen kan worden dan wordt het voordeel duidelijk!

**Oefening 47** *Stel dat wij 10.000.000.000 records willen opslaan die door een unieke sleutel als een getal voorgesteld kunnen worden. Een blok heeft 4096 bytes en wij willen de toppen zo kiezen dat ze in één blok passen.*

*Eerst werken wij met een B-tree waar wij de sleutels (5 bytes), de boompointers (8 bytes) en de pointers naar de records (8 bytes) in een top plaatsen.*

*Wat is de diepte van deze boom in het slechtste geval?*

*Dan werken wij met een B+-tree waar wij de pointers dus maar één keer nodig hebben – als het om interne toppen gaat, zijn het pointers naar toppen en anders naar records. Wat is de maximale diepte van deze boom?*



Figuur 14: Een voorbeeld voor het toevoegen in een B+-tree.

## 4 Algoritmen voor strings

Iedereen van jullie heeft zeker al algoritmen voor strings gebruikt – of het nu *query replace* in een editor is of het zoeken van woorden in een tekst. Maar ook op andere, minder duidelijke plaatsen zijn stringalgoritmen belangrijk – bv. in de bioinformatica waar de strings bv. het DNA voorstellen of in databanken waar de strings een identifier – bv. van een scheikundig molecuul kunnen zijn. Wat precies met de strings gedaan wordt – of getest moet worden of twee strings identiek zijn (misschien de gemakkelijkste taak), of of een string deelstring van een andere string is, of of één string een *mutatie* van een andere string kan zijn, etc. . . verschilt natuurlijk van geval tot geval. Wij zullen hier verschillende algoritmen zien – soms ook verschillende algoritmen voor dezelfde taak. De hoofdrede om deze algoritmen te leren is om de ideeën te leren kennen en ze achteraf zelfstandig misschien in andere omstandigheden te kunnen toepassen. Wij zullen hier – als een soort bewijs dat het nuttig is ideeën te kennen – bv. ook technieken zien die *in principe* geïnspireerd zijn door technieken die jullie al in DA1 en DA2 hebben gezien!

### 4.1 Exact string matching

In dit hoofdstuk zoeken wij een string  $z[]$  in een andere string  $t[]$ . Dit is dus de typische toepassing van zoekmachines of van een editor die een woord in een tekst zoekt. Wij zullen de string die gezocht wordt altijd  $z$  noemen ( $z$  voor zoeken) en de lengte van deze string  $m$  (de string is dus  $z[0]z[1]z[2] \dots z[m-1]$ ). De string waarin gezocht wordt (vaak ook de tekst genoemd) heet altijd  $t$  ( $t$  voor tekst) en de lengte wordt als  $n$  genoteerd (deze string is dus  $t[0]t[1] \dots t[n-1]$ ).

Natuurlijk heeft iedereen onmiddellijk een idee hoe je een string kan zoeken:

#### Algoritme 8 (*Brute kracht*)

```
strings_gelijk(z,t,start)
// test of de eerste |z| tekens van t -- beginnend met
// positie start -- gelijk zijn aan die van z
// het wordt gesteld dat in t na start nog ten minste even
// veel tekens staan als in s en dat |z| gekend is
{
for (i=0; i< |z|; i++) if (z[i]!=t[start+i]) return FALSE;

return TRUE;
```

```

}

is_contained(z,t)
// test of de string z in de string t voorkomt
{
for (i=0; i<= |t|-|z|; i++)
    if strings_gelijk(z,t,i) return TRUE;

return FALSE;

}

```

Hier worden de strings vanaf  $n - m + 1$  startposities vergeleken en elke vergelijking vraagt ten hoogste  $m$  stappen, de totale complexiteit is dus  $O((n - m + 1) * m)$ .

**Opmerking 3** *Algoritme 8 toegepast op een tekst van lengte  $n$  en een zoekstring van lengte  $m$  vraagt in het slechtste geval  $O((n - m + 1) * m)$  stappen.*

Als wij  $m$  als constant beschouwen (normaal zal de zoekstring heel klein zijn in vergelijking met de tekst) is dat  $O(n)$  en dus zeker niet slecht. Maar er zijn twee redenen om toch betere algoritmen te zoeken. Ten eerste is het wel een beetje gesjoemeld als wij  $m$  verwaarlozen – de zoekstring **kan** in principe ook  $\Theta(n)$  zijn – bv.  $n/2$  – en dan mag  $m$  natuurlijk niet als constant beschouwd worden en wordt de complexiteit in feite  $O(n^2)$ . En ten tweede is voor praktische doeleinden een grote constante factor voor de complexiteit ook niet echt goed en vooral in DA3 ligt de klemtoon toch sterk op de praktische kant van algoritmen!

**Oefening 48** *Als de binnenste lus van algoritme 8 echt  $m$  stappen nodig heeft, hebben wij het woord gevonden – dus zal het algoritme stoppen. En als de lus bijna  $m$  stappen doet, zouden wij verwachten dat misschien de volgende keren dat de functie wordt gebruikt al heel snel een tegenstrijdigheid gevonden wordt. Het is dus niet onmiddellijk duidelijk dat het algoritme echt soms  $\Theta(n^2)$  stappen nodig heeft. Daarom moet dat aangetoond worden: Geef een geparametriseerde reeks van voorbeelden waarvoor dit algoritme in feite  $\Theta(n^2)$  stappen nodig heeft. Natuurlijk moet je ook aantonen dat  $\Theta(n^2)$  stappen nodig zijn.*

**Oefening 49** • *Herschrijf Algoritme 8 zo dat de zoekstring ook één of meerdere wildcards ? mag bevatten. Wij stellen dat ? geen deel uitmaakt*

van de tekst en dat elk leetterteken erdoor gematcht wordt. Voor `?uizen` zou dus bv. `huizen` maar ook `buizen` een match in de tekst zijn.

*Wat is de complexiteit van jouw algoritme?*

- Geef een algoritme zodat de zoekstring ook één of meerdere wildcards `*` mag bevatten waarbij `*` alle deelstrings matcht – ook de lege string. Voor `ver*uizen` zouden dus bv. `verhuizen` en `vergeet de muizen` een match in de tekst zijn. Jouw algoritme moet een arbitraire string in de tekst teruggeven die de zoekstring matcht als die bestaat en anders de lege string.

*De complexiteit van jouw algoritme mag ten hoogste de complexiteit van Algoritme 8 zijn.*

In DA1 en DA2 werd gezegd dat één van de bedoelingen van deze lessen is om technieken en ideeën te zien die je achteraf in een andere context kan toepassen om snelle algoritmen te ontwikkelen. Het volgende algoritme is een voorbeeld daarvan.

## Het algoritme van Rabin-Karp

Als wij de kost van de functie `strings_gelijk()` in Algoritme 8 constant zouden kunnen maken of de functie alleen maar in heel zeldzame gevallen zouden moeten toepassen, zou de factor  $m$  in de complexiteit verdwijnen en wij zouden een complexiteit van  $O(n)$  hebben. Het algoritme van Rabin-Karp doet dat in feite niet **gegarandeerd** – de complexiteit van het **slechtste geval** blijft dezelfde – maar in de praktijk is de kost van de met de functie `strings_gelijk()` corresponderende functie wel constant. Misschien doet jullie dat al aan een techniek uit DA1 denken waar het slechtste geval ook geen verbetering ten opzichte van de standaardtechnieken was, maar die in de praktijk wel heel nuttig is: hashing. En het idee van hashing – dat ook al voor Bloom-filters toegepast werd – is wat wij ook hier zullen gebruiken.

Voordat wij de zoekstring en een deelstring van de tekst vergelijken, vergelijken wij eerst de waarde van een hashfunctie van deze strings. Als die niet gelijk is, zijn de strings zeker verschillend en moeten wij ze ook niet expliciet vergelijken. Dat betekent wel dat wij misschien een hashfunctie moeten berekenen voor alle deelstrings van lengte  $m$  in de tekst  $t$ . Als deze berekeningen even duur of duurder zouden zijn dan het vergelijken van de twee strings zou het natuurlijk geen voordeel zijn. Het trucje is dus dat wij om de hashfunctie voor de string die op positie  $i + 1$  begint te berekenen de hashwaarde voor de string die op positie  $i$  begint gewoon updaten.

Wij zullen als voorbeeld een hashfunctie voor het alfabet  $\Sigma = \{0, 1\}$  geven die al door Rabin en Karp werd gesuggereerd – wij werken dus met bitstrings. Maar het is duidelijk (zie oefeningen) dat het voor elk alfabet op een gelijkaardige manier gedaan kan worden.

Kies een natuurlijk getal  $q$  (bv. een priemgetal) en interpreteer een bitstring  $s$  (in ons geval kan dat ofwel  $z[]$  ofwel een deelstring van  $t[]$  zijn) gewoon als voorstelling van een getal  $g(s)$ . De hashwaarde  $h(s)$  van de string is dan gewoon  $h(s) = g(s) \bmod q$ . Als de string  $s = s_0s_1s_2 \dots s_{m-1}$  is dat als formule

$$h(s) = \left( \sum_{i=0}^{m-1} (s_i 2^{(m-1)-i}) \right) \bmod q$$

De som ziet er misschien ingewikkeld uit, maar het is alleen maar het getal dat je krijgt als je de bitstring interpreteert als de voorstelling van een binair getal.

Als wij dus de hashfunctie willen berekenen van de string van lengte  $m$  die op positie  $p$  in de tekst  $t_0t_1 \dots t_n$  begint, dan is dat

$$h(t_p \dots t_{p+m-1}) = \left( \sum_{i=0}^{m-1} (t_{p+i} 2^{(m-1)-i}) \right) \bmod q$$

Hierbij hebben wij de posities als index geschreven in plaats van als arrayindex (dus  $t_i$  in plaats van  $t[i]$ ) omdat dat formules beter leesbaar maakt.

Als wij gebruiken dat

$$(a + b) \bmod q = (a + (b \bmod q)) \bmod q$$

en

$$(a * b) \bmod q = (a * (b \bmod q)) \bmod q$$

dan krijgen wij voor de string die op positie  $p + 1$  begint:

$$\begin{aligned} h(t_{p+1} \dots t_{p+m}) &= \left( \sum_{i=0}^{m-1} (t_{p+1+i} 2^{(m-1)-i}) \right) \bmod q \\ &= \left( \sum_{i=0}^{m-2} (t_{p+1+i} 2^{(m-1)-i}) + t_{p+m} \right) \bmod q \\ &= (-t_p 2^m + \sum_{i=0}^{m-2} (t_{p+1+i} 2^{(m-1)-i}) + t_{p+m}) \bmod q \\ &= (-t_p 2^m + \sum_{i=0}^{m-1} (t_{p+i} 2^{m-i}) + t_{p+m}) \bmod q \\ &= (2 \sum_{i=0}^{m-1} (t_{p+i} 2^{(m-1)-i}) - t_p 2^m + t_{p+m}) \bmod q \\ &= ((2 \sum_{i=0}^{m-1} (t_{p+i} 2^{(m-1)-i}) \bmod q) - t_p 2^m + t_{p+m}) \bmod q \\ &= (2h(t_p \dots t_{p+m-1}) - t_p 2^m + t_{p+m}) \bmod q \end{aligned}$$

Deze berekening geeft ons een heel snelle en gemakkelijke manier om de hashwaarde  $h(t_{p+1} \dots t_{p+m})$  in constante tijd te berekenen als wij de hashwaarde  $h(t_p \dots t_{p+m-1})$  al kennen. Alleen voor het berekenen van de hashwaarde van



de zoekstring  $z$  en van  $h(t_0 \dots t_{m-1})$  hebben wij dus tijd  $O(m)$  nodig – alle andere hashwaarden die wij nodig hebben, kunnen in constante tijd berekend worden.

De volgende pseudocode beschrijft het algoritme van Rabin-Karp

**Algoritme 9** (*Rabin-Karp*)

```
strings_gelijk(z,t,start)
// test of de eerste |z| tekens van t -- beginnend met positie
// start -- gelijk zijn aan die van z
// het wordt gesteld dat in t na start nog ten minste even veel
// tekens staan als in z en dat |z| gekend is
{
for (i=0; i<|z|; i++) if (z[i]!=t[start+i]) return FALSE;

return TRUE;
}

is_contained_Rabin_Karp(z,t)
// test of de string s in de string t voorkomt
{
bereken h(z)

for (i=0; i<=|t|-|z|; i++)
{ bereken h(t_i...t_{i+m-1})
  // als i>0 kan dat gewoon door de oude waarde van h()
  // te updaten
  if (h(t_i...t_{i+m-1})==h(z))
    { if strings_gelijk(z,t,i) return TRUE; }
}
return FALSE;
}
```

De extra tijd in vergelijking met Algoritme 8 voor het berekenen van de hashfuncties is  $O(m)$  voor het berekenen van  $h(z)$ ,  $O(m)$  voor het berekenen van  $h(t_0 \dots t_{m-1})$  en dan nog  $O(n - m)$  voor het updaten van  $h()$  – samen dus  $O(n)$ . De functie `strings_gelijk()` vraagt nog altijd tijd  $O(m)$ , in het slechtste geval is de complexiteit dus nog altijd  $O((n - m + 1) * m)$ . Maar als wij ervan uitgaan dat de hashfunctie goed is en dus de waarden van de

deelstrings ongeveer gelijk zijn verdeeld over de  $q$  mogelijkheden, dan wordt de functie maar in  $(n - m + 1)/q$  gevallen toegepast. De totale tijd die ervoor nodig is, is dus  $O(m * (n - m + 1)/q)$  en als wij  $q \geq m$  kiezen dus  $O(n - m + 1)$ . De buitenste lus vraagt ook tijd  $O(n - m + 1)$  – alles samen vraagt als wij deze veronderstellingen doen dus een tijd van  $O(n)$ .

**Opmerking 4** *Het algoritme van Rabin Karp vraagt in het slechtste geval (als de waarden van de hashfunctie heel slecht verdeeld zijn) – net zoals het Algoritme 8 –  $O((n - m + 1) * m)$  stappen.*

*In de praktijk – en veronderstellend dat de waarden van de hashfunctie gelijkmatig verdeeld zijn en  $q$  voldoende groot is – vraagt het algoritme tijd  $O(n)$ .*

**Oefening 50** *Je zoekt de string  $z = 00101$  en past het algoritme van Rabin Karp toe met  $q = 13$ .*

*Geef een voorbeeld van een tekst  $t$  van lengte  $n$  waarvoor het algoritme bijzonder slecht presteert als je naar  $z$  zoekt.*

## Oefening 51

- *Beschrijf precies hoe je de hashwaarden berekent en update om te garanderen dat je geen problemen met overflow van integers krijgt.*
- *Het argument dat bewijst dat het updaten van de hashwaarde in constante tijd kan gebeuren telt in feite alleen het aantal integer operaties. Maar voor arbitrair grote waarden van  $q$  passen de hashwaarden niet noodzakelijk in een integer van beperkte grootte. Wat is de invloed op de complexiteit als je voor de hashwaarden de arbitrair grote integers door integers van beperkte grootte (bv. 64 bit) moet simuleren? Zal dat in de praktijk gebeuren? Waarom (of waarom niet)?*
- *Stel dat je met  $q$  werkt dat groter is dan de maximale som in de formule – dat komt dus neer op het schrappen van de modulo-berekening. Dan hoef je als de hashwaarden overeenstemmen de strings zelfs niet meer vergelijken omdat ze zeker gelijk zijn. Is dat een goed idee?*

**Oefening 52** *Stel dat je als alfabet niet alleen  $\Sigma = \{0, 1\}$  hebt maar een algemeen (eindig) alfabet  $\Sigma$ . Definieer voor  $\Sigma$  een hashfunctie die in constante tijd geüpdatet kan worden.*

*Denk daarbij erover na, of het belangrijk is dat in de som als basis voor de exponenten een 2 wordt genomen. Werken de argumenten ook met een andere basis – bv een ander priemgetal?*

## Oefening 53

- *Wijzig het algoritme van Rabin-Karp zo dat de zoekstring ook één wild-card ? mag bevatten. Wij stellen dat ? geen deel uitmaakt van de tekst en dat elk letterteken erdoor gematcht wordt (zoals in oefening 49). Ook hier moet het updaten van de hashfunctie (behalve in het begin) in constante tijd mogelijk zijn.*
- *Wijzig het algoritme van Rabin-Karp zo dat de zoekstring ook één of meerdere wildcards \* mag bevatten (zoals in oefening 49).*

*Jouw algoritme moet een arbitraire string in de tekst teruggeven die de zoekstring matcht als die bestaat en anders de lege string.*

*De complexiteit van jouw algoritme mag ten hoogste de complexiteit van het algoritme van Rabin-Karp zijn.*

**Oefening 54** *Wijzig het algoritme van Rabin-Karp zo dat je efficiënt naar een langste match kan zoeken.*

## Het algoritme van Knuth-Morris-Pratt

Als jullie naar het brute kracht algoritme 8 kijken dan komt de factor  $m$  van de functie `strings_gelijk()`. Maar deze functie is alleen maar duur als het deel in de tekst waarmee wij het vergelijken in het begin gelijk is aan de string die wij zoeken. In principe duikt het probleem dus altijd voor een deelstring op die zo begint als de zoekstring – die wij dus al op voorhand heel goed kennen. En dat betekent dat met een beetje preprocessing voor deze string zeker een voordeel te halen is. En dat is het principe van het door D. Knuth, J.H. Morris en V. Pratt in 1977 ontwikkelde algoritme.

Het principe wordt duidelijk als je naar een voorbeeld kijkt: Wij zoeken de string `zoek_me` in een tekst. Stel dat je probeert of die in de tekst  $t$  vanaf positie  $i$  voorkomt. Als al na 1 stap de vergelijking fout teruggeeft – dus  $t[i] \neq z$  moet je de volgende keer vanaf positie  $i + 1$  zoeken. Die heb je ten slotte nog niet gezien en het is dus mogelijk dat daar een match staat. Maar stel nu bv. dat de posities  $0, \dots, 4$  van de zoekstring nog overeenkomen en de vergelijking pas op positie 5 mislukt. Dat betekent dat vanaf positie  $i$  de tekst `zoek_` in de tekst staat maar niet de tekst `zoek_m`. Nu heb je de volgende posities wel al gezien en daaruit kan je een voordeel halen: je weet dat bv. tot positie  $t[i + 4]$  geen  $z$  voorkomt en je dus de string ten vroegste vanaf positie  $i + 5$  terug kan vinden. En je kan dat al op voorhand weten: **altijd** als ik `zoek_me` opzoek en de vergelijking op positie 5 voor de

eerste keer mislukt, kan ik de volgende 4 posities als mogelijke startposities overslaan.

### Het basisidee:

Ik wil op voorhand voor mijn string een verschuivingstabel  $V[]$  opstellen. Deze tabel bevat voor elke positie van de zoekstring het aantal vakjes waarmee ik mijn startpositie in de volgende iteratie zal mogen opschuiven als er de eerste mismatch optreedt op die positie in de zoekstring.

In het geval dat het eerste letterteken maar één keer in de zoekstring  $z[]$  opduikt (zoals in `zoek_me`) is de tabel  $V[]$  gemakkelijk te berekenen:  $V[0] = 1$  en  $V[j] = j$  voor  $0 < j \leq |s|$ . Je weet namelijk dat het eerste letterteken verschilt van de volgende  $j$  lettertekens – een vergelijking die vroeger in de tekst start dan  $j$  tekens later zou dus al voor het eerste letterteken van  $z[]$  een fout opleveren!

Voor andere teksten, als de zoekstring bijvoorbeeld `barbados` is, is het al iets moeilijker

We zullen nu eerst een precieze definitie van  $V[]$  ontwikkelen waarmee we verder kunnen werken.

Stel dat wij de posities  $t[i+0], \dots, t[i+j]$  al met  $z[0], \dots, z[j]$  vergeleken hebben en dat  $z[j] \neq t[i+j]$  de eerste mismatch is. Hierbij is  $j < |z|$  aangezien we anders het zoekwoord al volledig gematcht zouden hebben. Een verdere potentiële match zal dus zeker na positie  $j-1$  eindigen aangezien tot deze positie nog geen match gevonden werd. Nieuwe startposities die nog geen mismatch voor positie  $j-1$  geven, zijn dus posities  $s$  waarbij  $z[0] = z[s], \dots, z[j-1-s] = z[j-1]$ . Als zo'n positie  $s$  niet bestaat, m.a.w. er is geen interne match binnen de zoekstring tot positie  $j-1$ , dan is  $j$  de eerst mogelijke startpositie waarop een verdere match mogelijk is. Aangezien we de eerst mogelijke startpositie zoeken, nemen we de kleinst mogelijke  $s$  die aan de voorwaarden voldoet. Of formeel:

**Definitie 3** Voor een gegeven zoekstring  $z$  is de verschuivingstabel  $V[]$  gedefinieerd als

$V[0] = 1$  en voor  $j > 0$

$$V[j] = \min\{s \mid (0 < s < j \text{ en } z[0] = z[s], \dots, z[j-1-s] = z[j-1]) \\ \text{of } s = j\}$$

Als je op een positie  $i$  met  $0 < i < V[j]$  begint dan weet je dat er ten laatste op positie  $j-i$  een botsing is en dus de string vanuit deze positie niet gevonden kan worden.

Voor het voorbeeld `barbados` betekent dat

	b	a	r	b	a	d	o	s
positie j	0	1	2	3	4	5	6	7
V[j]	1	1	2	3	3	3	6	7

Wij weten nu hoeveel lettertekens je kan overslaan als je een match tot positie  $j - 1$  hebt maar geen tot positie  $j$  – maar er is nog een observatie die nodig is om ervoor te zorgen dat het algoritme echt lineair is:

**Observatie:** Stel dat er een match van  $z[]$  tot positie  $j - 1$  in  $t[]$  vanaf positie  $i$  staat. Dan kan een complete match ten vroegste op positie  $i + V[j]$  beginnen (dat hadden wij al). **Maar:** als  $V[j] < j$  dan weten wij van de posities  $i + V[j], \dots, i + j - 1$  al dat ze overeenkomen met  $z[V[j]], \dots, z[j - 1]$  – ze werden al vergeleken. Maar volgens de definitie van  $V[]$  komen die dan ook overeen met  $z[0], \dots, z[j - 1 - V[j]]$ . Als wij de string vanaf positie  $i + V[j]$  zoeken, moeten wij dus de eerste  $j - V[j]$  tekens niet opnieuw vergelijken. Wij moeten pas beginnen op de positie waar de eerste mismatch opdook.

Er kan dus op twee manieren bespaard worden: aan de ene kant wordt `vergelijk_strings()` minder vaak opgeroepen (door sommige startposities over te slaan) en aan de andere kant zijn de oproepen goedkoper omdat niet altijd vanaf het begin vergeleken moet worden.

Stel op dit moment dat wij de tabel  $V[]$  efficiënt kunnen berekenen. Dan is het algoritme van Knuth-Morris-Pratt als volgt:

#### Algoritme 10 (*Knuth-Morris-Pratt*)

```
vergelijk_strings(z,t,start,gelijke_tekens)
// zoekt of de string z in de tekst t vanaf positie start voorkomt.
// Het wordt verondersteld dat de tekens 0...(gelijke_tekens-1)
// in z en t (vanaf start) gelijk zijn.
// Het wordt bovendien verondersteld dat in t na start nog ten
// minste even veel tekens staan als in z en dat |z| gekend is.
// De functie geeft de index van de eerste mismatch terug
{
  for (i=gelijke_tekens; i<|z|; i++)
    if (z[i]!=t[start+i]) return i;

  return |z|;
}

is_contained_Knuth_Morris_Pratt(z,t)
```

```

// test of de string z in de string t voorkomt
{
  bereken_V(z)

  start=0;
  gelijke_tekens=0;

  while (start<=|t|-|z|)
  {
    eerste_mismatch=vergelijk_strings(z,t,start,gelijke_tekens);
    if (eerste_mismatch=|z|) return TRUE;

    start=start+V[eerste_mismatch];
    if (eerste_mismatch=0) gelijke_tekens=0; // speciaal geval
    else gelijke_tekens=eerste_mismatch-V[eerste_mismatch];
  }

  return FALSE;
}

```

Het geval waar de mismatch al op positie 0 opduikt is in zekere zin speciaal omdat daar ook de argumenten waarom wij kunnen opschuiven anders zijn. Het is het enige geval waar wij opschuiven **na** de positie van de eerste mismatch.

Het beste is natuurlijk de werking van het algoritme eens te illustreren aan de hand van een voorbeeld. Wij zoeken de string **barbados** in **barbaarse\_barbecue\_in\_barbados**. De eerste positie die vergeleken moet worden wordt met een v gemarkeerd.

```

v
barbaarse_barbecue_in_barbados
barbados
    5

```

De eerste mismatch is op positie 5. Omdat  $V[5] = 3$  kunnen wij met 3 posities opschuiven en moeten de eerste  $5 - 3 = 2$  tekens niet meer vergeleken worden:

```

v
barbaarse_barbecue_in_barbados
  barbados
    2

```

De eerste mismatch is nu op positie 2. Omdat  $V[2] = 2$  kunnen wij met 2 posities opschuiven en moeten de eerste  $2 - 2 = 0$  tekens niet meer vergeleken worden – dat helpt in dit geval dus niet.

```

      v
barbaarse_barbecue_in_barbados
      barbados
      0

```

Wij hebben onmiddellijk een mismatch, wij kunnen dus met  $V[0] = 1$  positie opschuiven en moeten 0 posities (speciaal geval) niet vergelijken:

```

      v
barbaarse_barbecue_in_barbados
      barbados
      0

```

Wij hebben opnieuw onmiddellijk een mismatch en kunnen dus met  $V[0] = 1$  positie opschuiven. Dat gaat door totdat wij de **b** van **barbecue** bereikt hebben:

```

      v
barbaarse_barbecue_in_barbados
      barbados
      4

```

De eerste mismatch is nu op positie 4. Omdat  $V[4] = 3$  kunnen wij met 3 posities opschuiven en moeten de eerste  $4 - 3 = 1$  tekens niet meer vergeleken worden:

```

      v
barbaarse_barbecue_in_barbados
      barbados
      1

```

De eerste mismatch is nu op positie 1. Nu zal dat zo doorgaan en er zal altijd onmiddellijk een mismatch gedetecteerd worden totdat de volgende **b** bereikt wordt – en dan wordt de match ook gevonden:

```

      v
barbaarse_barbecue_in_barbados
      barbados
      8

```

Maar wat is de complexiteit van dit algoritme? Op dit moment stellen wij dat het berekenen van de tabel in tijd  $O(m)$  kan gebeuren – dat zullen wij later natuurlijk nog bewijzen – en kijken wij alleen naar de routine die de string in de tekst echt zoekt.

**Stelling 5** *Het algoritme van Knuth, Morris en Pratt vraagt tijd  $O(n)$  als in een tekst  $t$  van lengte  $n$  een zoekstring  $z$  van lengte  $m \leq n$  gezocht wordt.*

**Bewijs:** Wij veronderstellen hier dat “`bereken_V()`” in tijd  $O(m)$  (dus ook  $O(n)$ ) kan gebeuren. Dat zullen wij later nog zien.

De lus in `is_contained_Knuth_Morris_Pratt()` wordt duidelijk ten hoogste  $n - m$  keer doorlopen en de kost in de lus is constant als wij de kost van `vergelijk_strings()` niet meerekenen. Wij moeten dus alleen nog bewijzen dat alle oproepen van `vergelijk_strings()` samen ten hoogste tijd  $O(n)$  vragen.

De complexiteit van `vergelijk_strings()` wordt bepaald door het aantal vergelijkingen die in de routine gebeuren en daar is een belangrijk verschil tussen succesvolle vergelijkingen en vergelijkingen die een mismatch opleveren: je kan hetzelfde teken uit de tekst meerdere keren *zonder succes* vergelijken, maar (dat zullen wij tonen) maar één keer succesvol. Wij tellen hier de vergelijkingen met en zonder succes dus apart. Omdat `vergelijk_strings()` ten hoogste  $n - m$  keer wordt opgeroepen en omdat in elke oproep ten hoogste een vergelijking zonder succes gebeurt (bij een mismatch wordt onmiddellijk een waarde teruggegeven) zijn ten hoogste  $n - m$ , dus  $O(n)$  vergelijkingen zonder succes.

Voor de andere oproepen is het nuttig de ontwikkeling te zien van de absolute positie in de tekst die met de string vergeleken wordt. In het algoritme wordt altijd de positie relatief met de startpositie bijgehouden, maar als je naar de absolute positie in de tekst kijkt dan hebben jullie misschien al in het voorbeeld gemerkt dat die nooit kleiner wordt. In feite wordt op het moment dat een teken één keer succesvol vergeleken wordt dit teken nooit meer met de zoekstring vergeleken – en omdat er maar  $n$  lettertekens zijn, volgt uit deze observatie natuurlijk dat ook deze bewerkingen samen maar tijd  $O(n)$  vragen.

Wij moeten dus alleen nog de observatie bewijzen:

Stel dat `vergelijk_strings()` wordt opgeroepen met de variabelen `start` en `gelijke_tekens`. Dan wordt vergeleken vanaf positie `start+gelijke_tekens`.



Wij moeten nu bewijzen dat als de volgende oproep met `start'` en `gelijke_tekens'` is dat dan

$$\text{start}' + \text{gelijke\_tekens}' \geq \text{start} + \text{eerste\_mismatch}$$

omdat `start+eerste_mismatch` net de eerste positie **na** de succesvolle vergelijkingen is en dat garandeert dat al succesvol vergeleken tekens niet nog eens getest worden.

Als `eerste_mismatch=0` dan zijn er geen succesvolle vergelijkingen en `gelijke_tekens` en `gelijke_tekens'` zijn beide 0 dus `start'=start+1` – de bewering klopt dus. Anders geldt

$$\begin{aligned} \text{start}' + \text{gelijke\_tekens}' &= \\ \text{start} + V[\text{eerste\_mismatch}] + \text{eerste\_mismatch} - V[\text{eerste\_mismatch}] &= \\ \text{start} + \text{eerste\_mismatch} \end{aligned}$$

En precies dat wouden wij bewijzen.

■

In typische toepassingen die jullie kennen, zijn de zoekstrings veel korter dan de tekst – zelfs de brute kracht manier om  $V[]$  te berekenen (die in  $O(m^2)$  of zelfs  $O(m^3)$  draait afhankelijk van hoe je het implementeert) zou dan geen probleem zijn. Maar voor toepassingen waar het verschil tussen  $m$  en  $n$  niet zo groot is, moet de preprocessing wel efficiënt zijn – en het best in tijd  $O(m)$  gebeuren.

Wij zoeken niet voor elke index  $i$  van  $V[]$  waar de vroegste startpositie is, maar voor elke mogelijke startpositie – beginnend met de kleinste, dus 1 – kijken wij hoe ver we raken zonder mismatch. Voor de nog niet ingevulde indexen  $i$  tot en met de eerste mismatch is de startpositie dan de vroegste startpositie, dus de juiste waarde voor  $V[i]$ . In principe zoeken wij dus alleen maar een match van  $z[]$  in zichzelf. Daarvoor kan je **bijna** hetzelfde algoritme toepassen als voor het zoeken van de tekst. Je zoekt de tekst in zichzelf – dus zijn  $z[]$  en  $t[]$  dezelfde string. Omdat je nu naar matches van verschillende lengten zoekt, moet je wel tot het einde doorgaan. Belangrijk daarbij is dat de waarden van  $V[]$  die je gebruikt altijd waarden zijn die je al op voorhand hebt berekend.  $V[0] = V[1] = 1$  is voor elk zoekwoord hetzelfde – daarmee kan je dus beginnen.

Om de pseudocode eenvoudiger te houden gaan we ervanuit dat op het einde van  $z[]$  nog een speciaal `end_of_line` teken volgt dat verschilt van alle tekens in de string en dat je in  $V[]$  één teken meer kan invullen dan de lengte van

$z[]$  (ook al gebruik je dat achteraf niet). Op deze manier moeten wij geen testen in `vergelijk_strings()` invullen die garanderen dat je niet over de grenzen van de array leest, wat de pseudocode minder leesbaar zou maken.

**Algoritme 11** (*Berekenen van  $V[]$* )

```
bereken_V(z)
// bereken de verschuivingstabel van de string z
{
  V[0]=1;
  V[1]=1;

  start=1; // voor 0 zou je gewoon dezelfde string terugvinden
  gelijke_tekens=0;

  while (start<|z|-1)
  {
    eerste_mismatch=vergelijk_strings(z,z,start,gelijke_tekens);

    if (eerste_mismatch=0) V[start+1]=start+1;
    else
      for (i=gelijke_tekens+1; i<=eerste_mismatch; i++)
        V[start+i]=start;
    // de andere tekens werden al ingevuld

    start=start+V[eerste_mismatch]; // deze waarde is
                                   // gegarandeerd al ingevuld

    if (eerste_mismatch=0) gelijke_tekens=0; // speciaal geval
    else gelijke_tekens=eerste_mismatch-V[eerste_mismatch];
  }
}
```

Het is duidelijk dat `bereken_V()`  $O(m)$  tijdsbegrensd is omdat het extra werk in vergelijking met `is_contained_Knuth_Morris_Pratt()` (het invullen van  $V[]$ ) zeker  $O(m)$  tijdsbegrensd is. De `while()` lus loopt wel tot het einde, maar in feite hebben wij bij de analyse van `is_contained_Knuth_Morris_Pratt()` ook gewoon gesteld dat het ten hoogste tot het einde loopt. En het deel dat

`bereken_V()` vervangt is hier het invullen van maar 2 getallen in een array – dat is dus zeker constant.

De werking van `bereken_V()` zie je het best opnieuw aan een voorbeeld. Hoe berekent deze functie de tabel voor `barbados`? Wij beginnen:

	b	a	r	b	a	d	o	s
positie j	0	1	2	3	4	5	6	7
V[j]	1	1						

1

`barbados`

`barbados`

0

Wij hebben onmiddellijk een mismatch (`eerste_mismatch=0`), kunnen dus  $V[1+1] = 2$  invullen en  $V[0] = 1$  opschuiven.

	b	a	r	b	a	d	o	s
positie j	0	1	2	3	4	5	6	7
V[j]	1	1	2					

2

`barbados`

`barbados`

0

Wij hebben opnieuw onmiddellijk een mismatch, kunnen dus  $V[2+1] = 3$  invullen en  $V[0] = 1$  opschuiven.

	b	a	r	b	a	d	o	s
positie j	0	1	2	3	4	5	6	7
V[j]	1	1	2	3				

3

`barbados`

`barbados`

2

Nu hebben wij de eerste mismatch op positie 2. Wij kunnen dus  $V[3+1] = 3$  en  $V[3+2] = 3$  invullen en  $V[2] = 2$  posities opschuiven.

	b	a	r	b	a	d	o	s
positie j	0	1	2	3	4	5	6	7
V[j]	1	1	2	3	3	3		

5

`barbados`

`barbados`

0

Nu krijgen wij nog twee keer onmiddellijk een mismatch en hebben dan

	b	a	r	b	a	d	o	s
positie j	0	1	2	3	4	5	6	7
V[j]	1	1	2	3	3	3	6	7

**Oefening 55** Pas het algoritme van Knuth-Morris-Pratt toe op de zoekstring  $z[] = \text{volvoert}$  en de tekst  $t[] = \text{volvo\_vol\_vogels\_volvoert\_kunststukje}$ . Bereken de verschuivingstabel en pas het zoekalgoritme toe. Werk voldoende tussenstappen uit om te laten zien hoe het algoritme werkt.

**Oefening 56** Pas het algoritme van Knuth-Morris-Pratt toe op de zoekstring  $z[] = \text{ananas}$  en de tekst  $t[] = \text{bananen\_in\_ra'ananas'centrum}$ . Bereken de verschuivingstabel en pas het zoekalgoritme toe. Werk voldoende tussenstappen uit om te zien hoe het algoritme werkt.

**Oefening 57** Wat moet aan de pseudocode gewijzigd worden om alle matches van de zoekstring te vinden?

Werkt het algoritme van Knuth-Morris-Pratt ook in tijd  $O(n)$  als in een tekst  $t[]$  van lengte  $n$  alle voorkomens van een zoekstring  $z[]$  van lengte  $m \leq n$  gezocht worden?

**Oefening 58** Stel dat  $m \gg n$ . Geef argumenten waarom het algoritme van Knuth-Morris-Pratt zoals beschreven dan niet  $O(n)$  tijdsbegrensd is en beschrijf precies hoe je het moet wijzigen dat het dat wel is.

**Oefening 59** Bewijs: als het algoritme zonder de extra “observatie” wordt geïmplementeerd (dat is hetzelfde alsof je `gelijke_tekens` altijd gelijk aan 0 zet) zijn er gevallen waarvoor het algoritme  $\Theta(n^2)$  is.

**Oefening 60** Het algoritme houdt alleen rekening met de informatie waar je een match hebt en niet met het letterteken op de mismatch – hoewel dat ook al gelezen werd. Zou het helpen ook daarmee rekening voor de verschuivingstabel te houden?

**Oefening 61** Wijzig het algoritme van Knuth-Morris-Pratt zo dat het voor twee strings  $z[]$  en  $t[]$  de positie en de lengte van een langste match van een prefix van  $z[]$  in  $t[]$  teruggeeft.

Gezocht is dus een match van  $z[0]z[1] \dots z[i]$  waarvoor  $i$  zo groot mogelijk is.

**Oefening 62** Stel dat je met de wildcard `?` wil werken en het algoritme van Knuth-Morris-Pratt gewoon wijzigt door een vergelijking waarin een `?` betrokken is altijd te aanvaarden. Toon aan dat het gewijzigde algoritme niet juist werkt.

**Oefening 63** *Wijzig het algoritme van Knuth-Morris-Pratt zo, dat het met wildcards \* in de zoekstring  $z[]$  kan werken, maar om het even hoeveel wildcards gebruikt worden nog altijd gegarandeerd in tijd  $O(n)$  draait met  $n$  de lengte van de tekst  $t[]$ .*

*Bewijs expliciet dat jouw algoritme de gegeven complexiteit heeft.*

## Het algoritme van Boyer-Moore-Horspool

Het algoritme van Boyer-Moore-Horspool is een vereenvoudiging van het algoritme van Boyer and Moore (1977). Het algoritme van Boyer and Moore in zijn oorspronkelijke vorm was *bijna* even efficiënt als het algoritme van Knuth-Morris-Pratt en met een kleine wijziging van Galil in 1979 zelfs even efficiënt (dat is:  $O(n + m)$  voor elke toepassing). Maar het bijzondere aan het algoritme van Boyer-Moore-Horspool is dat het heel eenvoudig en gemakkelijk toe te passen is. Dat maakt het in de praktijk – waar je ook rekening moet houden met de constanten die in de  $O()$ -notatie verwaarloosd zijn – bijzonder snel.

Voor ons is het hier interessant te zien dat een eenvoudig idee soms tot een heel goed algoritme kan leiden – ook al is de prestatie in het slechtste geval misschien niet zo goed als die van andere – ingewikkeldere – algoritmen.

Het volgende algoritme is bijna identiek aan het brute kracht algoritme. Het enige verschil is dat de strings niet van positie 0 tot  $|z| - 1$  maar van  $|z| - 1$  tot 0 worden vergeleken. Daardoor kan soms vroeger en soms later gestopt worden – een echt voordeel is dat niet.

```
strings_gelijk_invers(z,t,start)
// test of t[start]...t[start+|z|-1] = z[0]...z[|z|-1]
// hierbij wordt vanaf het einde vergeleken

{
for (i=|z|-1; i>=0; i--) if (z[i]!=t[start+i]) return FALSE;

return TRUE;
}

is_contained(z,t)
// test of de string z in de string t voorkomt
{
for (i=0; i<= |t|-|z|; i++)
```

```

    if strings_gelijk_invers(z,t,i) return TRUE;

return FALSE;

}

```

Stel nu dat je bijvoorbeeld de string **barbara** in **barbados\_zoekt\_barbara** zoekt. De eerste keer dat de lus in **is\_contained()** wordt doorlopen vergelijk je

```

barbados_zoekt_barbara
barbara

```

en het eerste teken dat je vergelijkt is het laatste van **barbara**. Je vergelijkt dus **a** met **o**. Daarbij zie je dat de vergelijking niet alleen een fout oplevert – het is zelfs zo dat er in de zoekstring **barbara** helemaal geen **o** zit ... Het is dus zinloos vanaf een startpositie te vergelijken waar de **o** nog in het deel zit dat vergeleken wordt – de volgende startpositie met kans op succes is dus de positie na de **o**.

```

barbados_zoekt_barbara
    barbara

```

Nu vergelijken wij eerst een **t** met de laatste **a** – en zijn in dezelfde situatie: er is geen **t** in **barbara**, wij kunnen dus opschuiven tot na de **t**.

```

barbados_zoekt_barbara
        barbara

```

Nu vergelijken wij een **r** met de **a** en ook al is dat een mismatch, kunnen wij ons argument hier niet meer toepassen: **barbara** bevat wel degelijk een **r** (zelfs twee). Hoe ver kunnen wij dus opschuiven? Wij weten dat er alleen maar een kans is als onze startpositie zo is dat er een **r** op de positie van de **r** in de tekst terechtkomt. Als je van achteraan telt dan zie je dat je ofwel 1 ofwel 4 posities moet opschuiven. Als je 4 posities opschuift, toets je niet of er een mogelijke match is als je maar 1 positie opschuift. Voor lettertekens uit de zoekstring kiezen wij dus altijd voor het kleinste aantal posities (maar wel ten minste 1 positie) dat je kan opschuiven om een match van het laatste teken te hebben – in ons voorbeeld dus 1:

```

barbados_zoekt_barbara
            barbara

```

Deze ideeën zullen wij nu gebruiken om een verschuivingstabel op te stellen. Een groot verschil met Knuth-Morris-Pratt is daarbij dat je de verschuiving opzoekt gebaseerd op een teken in de tekst – en niet gebaseerd op een positie in de zoekstring zoals bij Knuth-Morris-Pratt. Als wij een zoekstring  $z[]$  hebben en de tabel  $V_H$  noemen, dan kunnen wij voor een letterteken  $x$  definiëren:

$$V_H(x) = \begin{cases} |z| & \text{als } x \text{ niet in } z[] \text{ voorkomt} \\ & \text{of alleen op positie } |z| - 1 \\ p - 1 & \text{als } |z| - p \text{ de laatste positie voor } |z| - 1 \\ & \text{is, waar } x \text{ in } z[] \text{ voorkomt} \end{cases}$$

Deze tabel kan je met de volgende pseudocode gemakkelijk en snel berekenen:

```
vul_VH(z)

{
// eerst voor alle bytes (ASCII tekens) de waarde |z| invullen:
for (i=0; i<256; i++) VH[i]=|z|;

// dan voor tekens die in z[] voorkomen de waarden overschrijven
// de laatste keer dat een teken voorkomt bepaald de waarde:

for (i=0; i<=|z|-2; i++) VH[z[i]]=|z|-i-1;
}
```

Deze tabel kan duidelijk heel gemakkelijk en snel ingevuld worden. Ook in de  $O()$ -notatie is dat lineair in  $m = |z|$ , dus  $O(m)$ .

Voor **barbara** is dat

$x =$	<b>a</b>	<b>b</b>	<b>r</b>
$V_H[x]$	2	3	1

en  $V_H[x] = 7$  voor alle andere tekens.

Daarmee kunnen wij nu het algoritme van Boyer-Moore-Horspool beschrijven:

```

is_contained_BMH(z,t)

{
  vul_VH(z);

  for (i=0; i<= |t|-|z|; )
    { if strings_gelijk_invers(z,t,i) return TRUE;
      i= i+ VH[t[i+|z|-1]]};
    }

  return FALSE;
}

```

Omdat wij de werking al voor het zoeken van `barbara` in `barbados_zoekt_barbara` gezien hebben, zullen wij hier niet nog een voorbeeld uitwerken.

**Oefening 64** *Pas het algoritme van Boyer-Moore-Horspool toe op de zoekstring `z[] = ananas` en de tekst `t[] = de_spa_in_ra'ananas'`.*

*Bereken de verschuivingstabel en pas het zoekalgoritme toe. Werk voldoende tussenstappen uit om te zien hoe het algoritme werkt.*

Natuurlijk zien jullie onmiddellijk, dat hier – vergelijkbaar met Knuth-Morris-Pratt – informatie die je in `strings_gelijk_invers()` verwerft ook gebruikt zou kunnen worden om soms verder op te schuiven of vergelijkingen te besparen. Dan gaan jullie al in de richting van het (iets ingewikkeldere) algoritme van Boyer-Moore.

**Oefening 65** *Gegeven een zoekstring `z[]` met  $|z| = m$  en een tekst `t[]` met  $|t| = n$ .*

- *Geef voorbeelden van `z[]` en `t[]` (die voor alle  $n \geq m > 0$  werken) die voor het algoritme van Boyer-Moore-Horspool het beste geval voorstellen. Wat is de complexiteit van het algoritme van Boyer-Moore-Horspool in dit beste geval ( $\Theta()$ -notatie).*
- *Geef voorbeelden van `z[]` en `t[]` (die voor alle  $n \geq m > 0$  werken) die voor het algoritme van Boyer-Moore-Horspool het slechtste geval voorstellen. Wat is de complexiteit van het algoritme van Boyer-Moore-Horspool in dit slechtste geval ( $\Theta()$ -notatie).*



## Het shift-AND algoritme

Het algoritme van Knuth, Morris, and Pratt kan gezien worden als een optimalisatie van het brute kracht algoritme dat wij misschien als eerste zouden proberen. Het volgende algoritme van R. Baeza-Yates en G. Gonnet is om meerdere redenen interessant: aan de ene kant is al de aanzet een beetje verrassend en ligt de klemtoon hier sterk op de praktijk – waardoor het natuurlijk heel goed bij DA3 past – en aan de andere kant zouden de technieken voor jullie ook een motivatie kunnen zijn in andere toepassingen waar jullie zelf algoritmen moeten ontwerpen eens na te denken of een vertaling naar bitvectoren misschien een snel algoritme zou opleveren! De theoretische slechtste-geval complexiteit is niet alleen  $O(n*m)$  – net zoals het brute-kracht algoritme – maar in dit geval heb je zelfs  $\Omega(n*m)$  stappen in het beste geval nodig (als je de hele tekst moet doorlopen) terwijl het brute-kracht algoritme in dat geval maar tijd  $\Omega(n)$  vraagt. Maar als wij rekening houden met toepassingen waar  $m$  relatief klein is, wordt het niet alleen lineair (dat wordt het brute-kracht algoritme ook), maar presteert het zelfs bijzonder goed omdat de constante in de lineaire functie heel klein is! Dit algoritme is vooral een goede keuze als je alle matches zoekt en er zijn er veel in de tekst waarin je zoekt.

In sommige boeken kan je dit algoritme als shift-AND terugvinden en in anderen als shift-OR – de redenen daarvoor zullen jullie in de oefeningen zien.

**Definitie 4** Gegeven een zoekstring  $z[]$  van lengte  $m$  en een tekst  $t[]$  van lengte  $n$ .

Dan definiëren wij de  $m \times n$  matrix  $M$  als  $M[i][j] = 1$  als de prefix uit de tekens t.e.m. positie  $i$  van  $z[]$  met de deelstring die in  $t[]$  op positie  $j$  eindigt gematcht kan worden en anders 0. Of formeel:

$$M[i][j] = \begin{cases} 1 & \text{als } z[0] = t[j-i], \dots, z[i] = t[j] \\ 0 & \text{anders} \end{cases}$$

Hierbij hebben wij de rijen en kolommen nummers vanaf 0 gegeven in plaats van vanaf 1.

Voor het geval  $t[] = \text{ananas}$  en  $z[] = \text{ana}$  krijgen wij

	a	n	a	n	a	s
a	1	0	1	0	1	0
n	0	1	0	1	0	0
a	0	0	1	0	1	0

In de laatste rij staat dus een 1 op positie  $j$  ( $M[m-1][j] = 1$ ) als en slechts als een match van het hele zoekwoord  $z[]$  op positie  $j$  van de tekst eindigt.

Deze matrix heeft  $n * m$  elementen, die alle  $n$  te berekenen en in te vullen vraagt dus in principe tijd  $O(m * n)$ .

Maar hoe berekenen wij deze matrix?

Wij zullen een dynamisch programmeren algoritme gebruiken dat bovendien – door bitvectoren te gebruiken – parallel werkt: meerdere van de  $M[i][j]$  worden tegelijk met hulp van vroegere waarden berekend.

Daarvoor definiëren (en berekenen) wij eerst de binaire karakteristieke vector  $C_x$  voor elke letter  $x$  in het alfabet. Deze vector heeft lengte  $m$  en

$$C_x[i] = \begin{cases} 1 & \text{als } z[i] = x \\ 0 & \text{anders} \end{cases}$$

$C_x$  houdt dus bij waar het teken  $x$  in  $z[]$  staat. In ons voorbeeld  $z[] = \text{ana}$  krijgen wij  $C_a = (101)$ ,  $C_n = (010)$  en  $C_x = (000)$  voor alle anderen tekens  $x$  in het alfabet.

De eerste kolom is gemakkelijk:  $M[i][0] = 0$  voor  $1 \leq i < n$  omdat op positie 0 natuurlijk geen match van lengte groter dan 1 kan eindigen. En  $M[0][0] = 1$  als en slechts als  $z[0] = t[0]$

Stel nu dat wij kolom  $j$  al hebben en kolom  $j + 1$  willen berekenen.  $M[0][j + 1] = 1$  als en slechts als  $t[j + 1] = z[0]$  – of equivalent  $M[0][j + 1] = C_{t[j+1]}[0]$ . Voor  $i > 0$  hebben wij  $M[i][j + 1] = 1$  als en slechts als er een match van lengte  $i$  (de posities  $0, \dots, i - 1$ ) tot positie  $j$  is **en**  $t[j + 1] = z[i]$ . Of anders geschreven:  $M[i][j + 1] = M[i - 1][j] \ \& \ C_{t[j+1]}[i]$  – en nu zien jullie misschien waar de naam van dit algoritme vandaan komt: als de  $M[][j]$  en de karakteristieke vector  $C_{t[j+1]}$  bitvectoren zijn dan is  $M[][j + 1] = \text{shift}^1(M[][j]) \ \& \ C_{t[j+1]}$  als de shift-operatie  $\text{shift}^1$  bit  $i - 1$  op positie  $i$  schuift en het nieuwe bit 0 op 1 zet. Als pseudocode ziet dat er als volgt uit – waarbij wij altijd alleen maar één kolom van de matrix bijhouden – zodra wij de nieuwe kolom berekend hebben, hebben wij de oude ten slotte niet meer nodig:

**Algoritme 12** (*Het shift-AND algoritme*)

```
shiftAND(z,t)
{
// test of de string z in de string t voorkomt
// geeft TRUE terug als ja anders FALSE
// alle karakteristieke vectoren C[x][] en ook de vector kolom[]
// zijn bitvectoren -- dus bv. een unsigned long int als m<=64
```

```

bereken_karakteristieke_vectoren(); // dat is straightforward

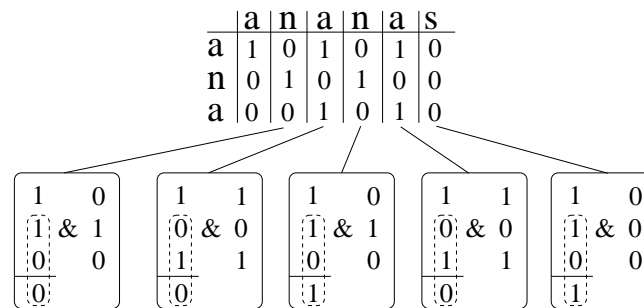
if (z[0]=t[0])
    { if (|z|=1) return TRUE;
      kolom=(1,0,0,...,0) // lengte |z|
    }
else kolom=(0,0,0,...,0)

for (j=1; j<|t|;j++)
{
    kolom[] = SHIFT^1(kolom[]) & C[t[j]][];
    if (kolom[|z|-1]=1) return TRUE;
}

return FALSE;
}

```

Hoe deze operaties tot de tabel voor ons voorbeeld leiden, zien jullie in Figuur 15.



Figuur 15: De stappen van het shift-AND algoritme toegepast op  $z[] = ana$  en  $t[] = ananas$  (waarbij na de eerste match niet gestopt wordt).

Als je deze operaties algemeen wil implementeren dan vraagt het berekenen van de volgende kolom – zoals verwacht – tijd  $O(m)$ . Als  $m$  naar boven door een constante begrensd is, is dat dus constant – maar ook dat is niet echt een voordeel in vergelijking met andere algoritmen. Maar als de lengte van de zoekstring ten hoogste het aantal bits van een woord in de computer is – en voor heel veel toepassingen mag je zeker stellen dat de zoekstring niet meer dan 64 tekens bevat – dan kan je de bitvectoren in één computerwoord voorstellen en de bitmanipulaties gebruiken ( $\ll, \gg, \&, |$ ) die de processor

kan uitvoeren. Dan zijn de bewerkingen niet alleen *constant* maar in feite *echt heel goedkoop*.

In gevallen waar je met heel lange woorden bezig bent, kan je eerst naar een prefix zoeken die nog in één computerwoord past. De lengte van een computerwoord is op dit moment meestal 64 bits. Als het om *gewone teksten* gaat, zal deze prefix vermoedelijk maar voor een heel klein deel van de mogelijke startposities gematcht kunnen worden. Voor de zeldzame gevallen waar die gematcht wordt, pas je voor de rest dan de brute kracht manier toe om te zien of de rest van de zoekstring hier aansluit.

Een andere mogelijkheid is met bitvectoren te werken die meer dan één computerwoord gebruiken.

**Oefening 66** *Geef een algoritme in pseudocode dat de verzameling van alle  $C_x$  in tijd  $O(m)$  kan berekenen waarbij  $m$  de lengte van het gezochte woord  $z[]$  is.*

**Oefening 67** *Schrijf de C-code die de shift-AND operaties implementeert voor een woordlengte op de computer van WORDSIZE bits.*

- voor een zoekwoord dat gegarandeerd ten hoogste WORDSIZE letters bevat
- voor een zoekwoord met onbeperkte lengte (maar werk hier wel met kolommen van dezelfde lengte als het woord!)

**Oefening 68** *Pas het shift-AND algoritme toe op het zoekwoord  $z[]$  =zoeken en de tekst  $t[]$  =iets\_zoets\_zoeken*

**Oefening 69** *Wijzig het shift-AND algoritme zo dat het ook met wildcards ? en \* kan werken en nog altijd even goed presteert.*

**Oefening 70** *Wij zijn eraan gewend dat 1 voor TRUE staat en 0 voor FALSE en zo is het shift-AND algoritme ook beter uit te leggen en te verstaan. Maar nu dat je het algoritme **hebt** verstaan: wissel in Definitie 4 de betekenis van 0 en 1.*

- Werk het algoritme uit voor deze gewijzigde definitie en geef de pseudocode.
- Schrijf de belangrijkste lijn van de C-code voor het algoritme dat het algoritme volgens deze nieuwe definitie implementeert.
- Welke versie zou je implementeren? Die met de interpretatie van 1 zoals in de tekst of zoals in de oefening? Waarom?

**Oefening 71** *Wijzig Algoritme 12 zo dat de lengte van een langste match beginnend met  $z[0]$  terug wordt gegeven.*

### shift-AND voor matches met fouten

Maar het shift-AND algoritme laat het ook toe matches nog meer te veralgemenen dan met wildcards: matches met een beperkt aantal mismatches – zonder op voorhand te specificeren waar die moeten zitten. Wij zullen hier het algoritme zien dat het programma *agrep* te gronde ligt en dat door S. Wu en U. Manber in 1992 ontwikkeld werd. Ook hier ligt de klemtoon op de praktijk – dus het geval waar het toegelaten aantal fouten begrensd en klein is.

Wij zullen het hier op dit moment alleen maar hebben over *mismatches* – dus waar een letterteken met een **ander** letterteken gematcht moet worden. Je kan het ook over *fouten* in de betekenis van ontbrekende lettertekens (in de tekst of de zoekstring) hebben.

**Definitie 5** *Gegeven een zoekstring  $z[]$  van lengte  $m$ , een tekst  $t[]$  van lengte  $n$  en een (klein) getal  $k$ .*

*Dan definiëren wij de  $m \times n$  matrix  $M^k$  als  $M^k[i][j] = 1$  als de prefix uit de tekens t.e.m. positie  $i$  van  $z[]$  met de deelstring die in  $t[]$  op positie  $j$  eindigt gematcht kan worden waarbij er ten hoogste  $k$  mismatches optreden en anders 0. Of formeel (waarbij wij gebruiken dat ten hoogste  $k$  mismatches betekent dat er ten minste  $i + 1 - k$  goede matches zijn):*

$$M^k[i][j] = \begin{cases} 1 & \text{als er indices } 0 \leq l_0 < l_1 < \dots < l_{i-k} \leq i \text{ bestaan zo dat} \\ & z[l_0] = t[j - i + l_0], \dots, z[l_{i-k}] = t[j - i + l_{i-k}] \\ 0 & \text{anders} \end{cases}$$

*Ook hierbij hebben wij de rijen en kolommen nummers vanaf 0 gegeven in plaats van vanaf 1.*

De matrix  $M^0[][]$  is dus de matrix  $M[][]$  die wij al kennen – hier laten wij 0 fouten toe en wij weten hoe wij die moeten berekenen.

In de laatste rij van  $M^k[][]$  staat dus een 1 op positie  $j$  ( $M^k[m-1][j] = 1$ ) als en slechts als een match met ten hoogste  $k$  fouten van het hele zoekwoord  $z[]$  op positie  $j$  van de tekst eindigt. Het idee lijkt dus heel sterk op het gewone shift-AND algoritme en de vraag is alleen maar hoe je  $M^k[][]$  het best berekent. De eerste kolom is gemakkelijk voor alle  $k > 0$ :  $M^k[][0] = (1, 0, \dots, 0)$

Er zijn 2 mogelijkheden waarom voor een  $0 < l \leq k$  de waarde van  $M^l[i][j] = 1$  is (wat betekent: er zijn ten hoogste  $l$  mismatches van een prefix van lengte

$i + 1$  met een string die op positie  $j$  eindigt). Wij veronderstellen dat  $j > 0$  en  $l > 0$  omdat wij voor deze waarden al weten hoe we  $M^l[i][j]$  berekenen. Bovendien kunnen wij ook stellen dat  $i > 0$  omdat  $M^l[0][j] = 1$  voor alle  $l > 0$ .

- als er ten hoogste  $l - 1$  mismatches van de prefix met lengte  $i$  tot positie  $j - 1$  zijn (dus  $M^{l-1}[i-1][j-1] = 1$ ) dan is het zonder belang of het letterteken op positie  $j$  gematcht wordt – er zijn zeker ten hoogste  $l$  fouten tot positie  $j$ .
- als er precies  $l$  mismatches van de prefix met lengte  $i$  tot positie  $j - 1$  zijn (dan is  $M^l[i-1][j-1] = 1$ ) en  $t[j] = z[i]$  dan zijn er ten hoogste  $l$  mismatches tot positie  $j$ .

In andere gevallen is er duidelijk geen match met de vereiste eigenschappen en moet  $M^l[i][j] = 0$  zijn. Wij moeten de matrices dus in een volgorde berekenen zodat de waarden voor  $l - 1$  en  $j - 1$  al bekend zijn.

De kolom  $M^l[.][j]$  kunnen wij dus berekenen als

$$M^l[i][j] = (1, 0, \dots, 0) \mid \text{shift}(M^{l-1}[i][j-1]) \mid (\text{shift}(M^l[i][j-1]) \ \& \ C[t[j]])$$

Hierbij kan  $\text{shift}()$  bv. dezelfde operatie zijn als  $\text{shift}^1()$  of het kan een  $\text{shift}$  zijn die de 0-bit op positie 0 schuift. Door de operatie  $(1, 0, \dots, 0) \mid$  wordt deze bit toch al op 1 gezet.

**Oefening 72** *Werk ook dit algoritme uit met de rol van 0 en 1 in Definitie 5 gewijzigd.*

**Oefening 73** *Schrijf de pseudocode voor dit algoritme op.*

**Oefening 74** *Pas het algoritme toe op  $z = \text{dank}$  en  $t = \text{een\_dansfeest}$  met 1 toegelaten mismatch.*

**Oefening 75** *Wij kunnen de definitie van toegelaten fouten in Definitie 5 ook zo beschrijven, dat er een perfecte match is met een woord dat wij uit het zoekwoord kunnen bouwen als wij ten hoogste  $k$  lettertekens door een ander letterteken vervangen. Stel dat het nu ook toegelaten is een letterteken uit de tekst te verwijderen (deletion) en dat dat ook als 1 fout telt (als de tekst dus bv. `ik_ben_peter` is dan zou dit algoritme voor de zoekstring `peer` zeggen dat er een match met 1 fout gevonden is). Geef de **exacte** definitie van  $\bar{M}^k[i][j]$  voor dit concept en geef het algoritme om  $\bar{M}^k[i][j]$  te berekenen.*

## 4.2 Benaderend string matching

Wij hebben al verschillende stappen in de richting van benaderend string-matching gedaan door wildcards of mismatches toe te laten. In dit deel willen wij het **eerst** over een speciaal geval hebben: het geval waar de rol van  $z[]$  en  $t[]$  symmetrisch is. Wij hebben hier dus niet één string die typisch kort is en een andere die lang is, maar ze zijn beide van ongeveer dezelfde lengte en wij willen weten *hoe verschillend* ze zijn. Je hebt bv. twee bestanden waarbij het ene bestand een gewijzigde versie van het andere bestand is en wilt weten hoe sterk het bestand gewijzigd werd. Dat wordt bv. door het Unix-commando *diff* berekend (waarbij *diff* ook nog zegt wat de wijzigingen zijn).

Als wij zeggen *hoe verschillend* dan hebben wij het in feite over een maat – of een afstand. Hoe je de definitie van een afstand kiest, hangt er sterk van af wat de context is. In de biologie zijn de strings bv. vaak DNA structuren en de afstand hangt ervan af hoeveel mutaties nodig zijn om van het ene DNA naar het andere te gaan. Maar je kan er ook best rekening mee houden hoe waarschijnlijk de mutaties zijn. Dat geldt natuurlijk ook voor teksten: je kan tellen hoeveel *wijzigingen* nodig zijn om de ene tekst in de andere te veranderen, maar dan is nog altijd niet duidelijk welke soort *wijzigingen* bedoeld zijn. Als de afstand moet beschrijven hoe waarschijnlijk het is dat het woord  $t[]$  in feite hetzelfde woord zou moeten zijn als  $t'[]$  dan zouden toegelaten wijzigingen bv. het weglaten van een letterteken, het toevoegen van een letterteken of het verwisselen van een letterteken zijn. Maar je zou ook hier waarschijnlijkheden kunnen gebruiken. Misschien is het waarschijnlijker dat lettertekens verwisseld worden die naast elkaar op het toetsenbord staan dan lettertekens die niet naast elkaar staan – of is het waarschijnlijker dat de volgorde van twee lettertekens fout is (bv. *volgorde* in plaats van *volgorde*) dan dat twee lettertekens verwisseld zijn (bv. *vilgorde* in plaats van *volgorde*), etc.

De afstand die wij hier gebruiken is dus gewoon één van vele mogelijke afstanden die allemaal in zekere omstandigheden zinvol kunnen zijn.

**Definitie 6** *Voor de editeerafstand zijn 3 bewerkingen toegelaten:*

- *Het vervangen van een letterteken door een ander letterteken. Deze bewerking hangt nauw samen met de mismatches die wij in samenhang met het shift-AND algoritme al besproken hebben.*
- *Het verwijderen van een letterteken. Als het verwijderde letterteken positie  $i$  had, houden alle tekens die voor het verwijderen een positie*

$j < i$  hadden de positie en alle tekens die een positie  $j > i$  hadden, hebben achteraf positie  $j - 1$ .

- *Het toevoegen van een letterteken. Als het nieuwe letterteken na het toevoegen positie  $i$  heeft, houden alle tekens die voor het toevoegen een positie  $j < i$  hadden de positie en alle tekens die een positie  $j \geq i$  hadden, hebben achteraf positie  $j + 1$ .*

De editeerafstand  $d(t, t')$  tussen twee strings  $t, t'$  is het kleinst mogelijke aantal bewerkingen van deze soort toegepast op  $t$  zodat het resultaat gelijk aan  $t'$  is.

Deze afstand wordt soms ook de Levenshtein afstand genoemd.

Wij zullen dit aan een voorbeeld illustreren. Wij gebruiken de strings  $t = \text{lessen}$  en  $t' = \text{feesten}$

l e s s e n	(1) vervang l door f
f e s s e n	(2) vervang s op positie 3 door t
0 1 2 3 4 5	
f e s t e n	(3) voeg nieuwe e op positie 1 toe
0 1 2 3 4 5	
f e e s t e n	
0 1 2 3 4 5 6	

Er bestaat dus een reeks met 3 stappen van  $t = \text{lessen}$  naar  $t' = \text{feesten}$ . En omdat er geen kortere reeks bestaat (gemakkelijk om te zien), geldt  $d(t, t') = 3$ .

**Opmerking 6** Voor strings  $t, t', t''$  geldt:

- (i)  $d(t, t') \geq 0$
- (ii)  $d(t, t') = 0$  als en slechts als  $t = t'$
- (iii)  $d(t, t') = d(t', t)$
- (iv)  $d(t, t'') \leq d(t', t) + d(t', t'')$
- (v)  $d(t, t') \leq \max\{|t|, |t'|\}$

Eigenschap (v) zegt dat de editeerafstand goed gedefinieerd is (hij kan niet oneindig zijn) en de eigenschappen (i) t.e.m. (iv) rechtvaardigen het over een *afstand* of een *metriek* te spreken.



**Bewijs:** De meeste punten volgen onmiddellijk uit de definitie. Het enige punt waarover wij misschien een beetje meer moeten nadenken is (iii) en dat punt zullen wij hier expliciet bewijzen:

Als  $d(t, t') = 0$  dan is  $t = t'$  en dan geldt zeker dat  $d(t, t') = d(t', t)$ . Maar je kan elke bewerking  $b$  zodat  $b$  toegepast op  $t$  de string  $t'$  oplevert omkeren: als de bewerking is dat  $t[i]$  vervangen wordt door letterteken  $x$  dan is  $t'[i] = x$  en kunnen wij als omgekeerde bewerking  $t'[i]$  vervangen door  $t[i]$ . Als een nieuw teken op positie  $i$  in  $t$  wordt toegevoegd dan kunnen wij als omgekeerde bewerking teken  $t'[i]$  verwijderen. Als teken  $t[i]$  wordt verwijderd dan kunnen wij omgekeerd in  $t'$  het teken  $t[i]$  op positie  $i$  toevoegen.

Stel dus dat  $d(t, t') = n \geq 1$ . Dan is er een reeks  $t = t_0, t_1, \dots, t_{n-1}, t_n = t'$  met  $d(t_i, t_{i+1}) = 1$  voor  $0 \leq i < n$ . Als wij de bewerkingen tussen twee opeenvolgende strings omkeren hebben wij een reeks die  $t'$  in  $t$  omvormt. Dus  $d(t', t) \leq n$ . Maar als  $d(t', t) < n$  zou op dezelfde manier volgen  $d(t, t') < n$  – een tegenstrijdigheid. Dus geldt  $d(t', t) = d(t, t')$ . ■

Om de editeerafstand efficiënt te berekenen hebben wij nog een lemma nodig:

**Lemma 7** *Gegeven  $t$  en  $t'$ . Dan is er een reeks van  $k = d(t, t')$  bewerkingen  $b_1, \dots, b_k$  die  $t$  naar  $t'$  wijzigt zodat als  $i_j$  de bij bewerking  $b_j$  betrokken index is, geldt  $i_1 \leq i_2 \leq \dots \leq i_k$ .*

Dit betekent dat wij beginnend aan de linkerkant van de string de string kunnen wijzigen en dan doorgaan naar de rechterkant waarbij wij nooit reeds vroeger bezochte delen opnieuw moeten wijzigen. Dat is natuurlijk heel nuttig voor een algoritme! Het betekent dat wij in feite de prefixen in volgorde van hun lengte wijzigen!

**Bewijs:** Dit kan je bewijzen door gewoon naar alle combinaties van bewerkingen te kijken waar eerst index  $i$  betrokken is en daarna index  $j$  waarbij  $j < i$ . In elk geval kan je ofwel bewijzen dat er een kortere reeks van bewerkingen is – wat een tegenstrijdigheid zou zijn – (bv. als je eerst een teken op positie  $i$  toevoegt en het dan door een ander teken vervangt) of dat je de bewerkingen kan vervangen door 2 bewerkingen met indices in volgorde.

Maar dit bewijs bevat noch interessante ideeën noch is het moeilijk – wij zullen het hier dus niet geven. ■

Eén van de doelen van de lessen over algoritmen is ook in te kunnen schatten welke technieken het best toegepast kunnen worden om een probleem op te lossen. Soms heb je dan zo'n gevoel dat het ene probleem op een zekere manier op het andere lijkt en dus misschien dezelfde technieken toegepast kunnen worden. Misschien hebben sommigen van jullie de indruk dat dit probleem een beetje op het probleem lijkt waar wij volgorden van matrices moesten bepalen om zo efficiënt mogelijk te kunnen vermenigvuldigen. Beide hebben een gelijkaardige *lineaire* structuur...

Wij zullen – net zoals voor het matrixvermenigvuldigingsprobleem en het shift-AND algoritme – ook hier dynamisch programmeren toepassen. Wij berekenen in feite niet alleen de editeerafstand tussen de strings  $t[]$  en  $t'[]$ , maar tussen alle prefixen van  $t[]$  en  $t'[]$ . Hier wordt dan heel duidelijk hoe dicht dit algoritme bij het shift-AND algoritme staat. Het is alleen jammer dat hier niet alleen de waarden 0 en 1 kunnen voorkomen die het mogelijk gemaakt hebben de bewerkingen door bitvectoren te paralleliseren. Als  $n = |t|$  en  $m = |t'|$  dan schrijven wij voor  $0 \leq i \leq n$  en  $0 \leq j \leq m$  de notatie  $d[i][j]$  voor de afstand tussen de prefix van lengte  $i$  van  $t[]$  (dus  $t[0], t[1], \dots, t[i-1]$ ) en de prefix van lengte  $j$  van  $t'[]$  (dus  $t'[0], t'[1], \dots, t'[j-1]$ ). De notatie duidt er al op dat wij dat in een 2-dimensionale array zullen bijhouden. De waarden  $d[i][0] = i$  en  $d[0][j] = j$  voor alle  $i, j$  zijn duidelijk omdat als één van de strings leeg is alle lettertekens ofwel verwijderd moeten worden ofwel toegevoegd moeten worden. Als wij veronderstellen dat wij op het moment dat wij  $d[i][j]$  willen berekenen de waarden van  $d[i'][j']$  voor alle tweetallen  $(i', j')$  die lexicografisch kleiner zijn dan  $(i, j)$  al kennen, dan kunnen wij als volgt argumenteren:

- als  $t[i-1] = t'[j-1]$  dan kan het laatste letterteken ongewijzigd blijven en dus  $d[i][j] \leq d[i-1][j-1]$ . Zie Oefening 77.

Als  $t[i-1] \neq t'[j-1]$  dan zijn er 3 bewerkingen die de laatste bewerking in een kortst mogelijke reeks kunnen zijn. Wij gebruiken hier de notatie  $t[k]$  om het letterteken te beschrijven dat voor de wijzigingen op positie  $k$  staat, maar merk op dat het voor de laatste wijziging door toevoegingen of verwijderingen op een helemaal andere positie kan staan.

- $t[i-1]$  vervangen door  $t'[j-1]$ . Dan is  $d[i][j] = d[i-1][j-1] + 1$ .
- $t'[j-1]$  aan de (misschien gewijzigde)  $t[0], t[1], \dots, t[i-1]$  toevoegen. Dan was de gewijzigde string ervoor  $t'[0], t'[1], \dots, t'[j-2]$  en die werd bereikt door  $t[0], t[1], \dots, t[i-1]$  te wijzigen. Wij hebben dus  $d[i][j] = d[i][j-1] + 1$ .

- $t[i-1]$  verwijderen. Dan wordt  $t[0], t[1], \dots, t[i-2]$  veranderd naar  $t'[0], t'[1], \dots, t'[j-1]$  en achteraf het overbodige letterteken  $t[i-1]$  verwijderd. Dus geldt  $d[i][j] = d[i-1][j] + 1$ .

Deze observaties kunnen nu gebruikt worden om een algoritme in pseudocode te formuleren. Daarbij zijn sommige dingen voor de betere verstaanbaarheid anders opgeschreven dan in een efficiënt programma – soms weet je bv. al dat een waarde beter is dan de waarde van **best** in de code.

**Algoritme 13** (*Dynamisch programmeren om de editeerafstand te berekenen*)

```
editeerafstand(t[], t'[])
{
// het wordt verondersteld dat t, t' beide niet leeg zijn en dat de
// lengte |t|, resp. |t'| is. De editeerafstand tussen t en t' wordt
// teruggegeven.

for (i=0; i<=|t|; i++) d[i][0]=i;
for (j=0; j<=|t'|; j++) d[0][j]=j;

for (i=1; i<=|t|; i++)
{
    for (j=1; j<=|t'|; j++)
    { // als k lettertekens beschouwd worden heeft het laatste
      // index k-1
      if (t[i-1]==t'[j-1]) buffer=d[i-1][j-1];
      else buffer=d[i-1][j-1]+1;
      d[i][j] = min{buffer, d[i-1][j]+1, d[i][j-1]+1};
    }
}

return d[|t|][|t'|];
}
```

Als wij de werking van dit algoritme eens beschouwen aan de hand van ons voorbeeld  $t[] = \text{lessen}$  en  $t'[] = \text{feesten}$  dan krijgen wij de volgende tabel. De lijn toont aan hoe de waarden verkregen zijn en ook één van de minst kostelijke mogelijkheden om de string **lessen** te wijzigen om de string **feesten** te verkrijgen.

		f	e	e	s	t	e	n
	0	1	2	3	4	5	6	7
l	1	1	2	3	4	5	6	7
e	2	2	1	2	3	4	5	6
s	3	3	2	2	2	3	4	5
s	4	4	3	3	2	3	4	5
e	5	5	4	3	3	3	3	4
n	6	6	5	4	4	4	4	3

Omdat de tijd voor één doorloop van de binnenste lus duidelijk constant is, hebben wij het volgende resultaat:

**Lemma 8** *De editeerafstand tussen twee strings van lengte  $n$  en  $m$  kan in tijd  $O(n * m)$  berekend worden.*

Voor de echte implementatie is het natuurlijk belangrijk geheugen te besparen en misschien niet de hele matrix in het geheugen te houden. Los zeker Oefening 78 op!

Het is niet alleen interessant te weten wat de editafstand is, maar ook de bewerkingen te kennen om de ene string naar de andere te wijzigen. Als je bv. een heel groot bestand hebt die op de twee computers zit dan is het misschien beter als die op de ene computer gewijzigd wordt alleen de nodige bewerkingen naar de andere computer te sturen dan het hele bestand.

De bewerkingen moeten ook niet altijd alleen met lettertekens gebeuren. In het Unix-commando *diff* worden bv. lijnen als characters geïnterpreteerd en de bewerkingen zijn dan verwijderen, toevoegen of vervangen van lijnen. . . Er zijn dus talrijke variaties en toepassingen van deze algoritmen. . .

**Oefening 76** *Gebruik het beschreven algoritme om de editeerafstand tussen  $t=roestig$  en  $t'=oesters$  te berekenen. Toon de tabel en duidt aan hoe de getallen in de tabel werden berekend. Beschrijf ook de mogelijke kortste reeksen van bewerkingen.*

**Oefening 77** *Stel dat  $t[] = t[0], t[1], \dots, t[i]$ ,  $t'[] = t'[0], t'[1], \dots, t'[j]$  en  $t[i] = t'[j]$ . Schrijf  $t_{i-1}$  voor  $t[0], t[1], \dots, t[i-1]$  en analoog  $t'_{j-1}$  voor  $t'[0], t'[1], \dots, t'[j-1]$ . Is dan altijd  $d(t, t') = d(t_{i-1}, t'_{j-1})$ ? Bewijs dat of geef een tegenvoorbeeld.*

**Oefening 78** *Herschrijf de pseudocode van Algoritme 13. Stel daarbij dat het antwoord op de vraag in Oefening 77 ja is en let erop zo weinig mogelijk geheugen te gebruiken. De strings kunnen natuurlijk te lang zijn om de hele matrix  $d[][]$  efficiënt in het geheugen te houden!*

**Oefening 79** *Bewijs gedetailleerd dat voor de volgende gevallen van twee opeenvolgende bewerkingen  $b_i, b_{i+1}$  die vervangen kunnen worden door minder bewerkingen of twee bewerkingen waarvan de tweede een teken met een groter index gebruikt dan de eerste.*

- $b_i$  verwijdert een teken op positie  $j$ ,  $b_{i+1}$  voegt een teken op positie  $j$  toe.
- $b_i$  voegt een teken op positie  $j$  toe,  $b_{i+1}$  verwijdert een teken op positie  $j' < j$ .
- $b_i$  verwijdert een teken op positie  $j$ ,  $b_{i+1}$  vervangt een teken op positie  $j' < j$  door een ander teken.

*Bepaal precies de paren van bewerkingen die in een reeks van minimale lengte zoals in Lemma 7 – dus met stijgende indices – gelijke indices moeten hebben en niet door even veel of minder bewerkingen met strict stijgende indices vervangen kunnen worden. Hier is een redenering die niet alle details geeft (dus de indices, etc) voldoende.*

**Oefening 80** *Definieer een nieuwe afstand  $d'()$  door behalve de drie al gekende bewerkingen ook nog een nieuwe bewerking toe te laten: het verwisselen van twee op elkaar volgende lettertekens. Met  $t=luil$  en  $t'=liul$  is dus  $d(t, t') = 2$  en  $d'(t, t') = 1$ . Beschrijf een algoritme om  $d'()$  te berekenen in pseudocode en geef uitleg waarom het algoritme juist is.*

**Oefening 81** *Wijzig de pseudocode van het algoritme zo dat ook de nodige bewerkingen om  $t[]$  naar  $t'[]$  te wijzigen worden bijgehouden.*

**Oefening 82** *Gegeven 2 woorden  $z_1[]$  en  $z_2[]$ . Het doel is deze twee woorden zo te wijzigen dat ze gelijk zijn. Maar de enige bewerking is dat je lettertekens tot  $z_1[]$  en  $z_2[]$  mag toevoegen. Beschrijf een efficiënt algoritme dat de minimale lengte van een woord  $z'[]$  bepaalt zodat  $z_1[]$  en  $z_2[]$  omgevormd kunnen worden naar  $z'[]$ .*

### De beste benaderende match

Nu zullen wij het terug over het probleem hebben matches van een korte string in teksten te vinden – alleen dat wij het hier over de beste benaderende match hebben en geen exacte matches. Hoe wij *benaderende matches* in de betekenis van *een gegeven aantal mismatches* efficiënt kunnen berekenen, hebben wij al gezien. Hier zullen wij het er nu over hebben hoe je voor een gegeven string  $z[]$  en een gegeven tekst  $t[]$  de deelstring van  $t[]$  met de kleinste editeerafstand kan vinden.

Het eerste idee zou misschien zijn het algoritme voor de editeerafstand tussen twee strings op alle deelstrings van  $t[]$  toe te passen. Maar met  $m = |z|$  en  $n = |t|$  zijn er  $(\sum_{i=0}^{n-1} (n-i)) + 1 = n(n+1)/2 + 1$  deelstrings van  $t[]$  (op positie  $i$  starten  $n-i$  mogelijke niet lege strings en dan nog de ene lege string 1 keer geteld). Om  $z[]$  met een string van lengte  $l$  te vergelijken hebben wij  $O(m * l)$  stappen nodig en als wij dat voor alle mogelijke lengten  $l$  opsommen, krijgen wij een totale kost van  $O(m * n^3)$ . En – veronderstellend dat het een goede bovengrens is – betekent dat natuurlijk dat een hierop gebaseerd programma vrij traag zal zijn voor grote teksten.

**Oefening 83** *Als je echt naar alle deelstrings van  $t[]$  kijkt, is dat natuurlijk niet bijzonder slim. Je ziet onmiddellijk dat je naar heel lange deelstrings niet moet kijken (waarom?). Maar wat is “heel lang”?*

*Geef een functie  $f()$  zodat voor alle  $z[]$  en  $t[]$  een deelstring in  $t[]$  bestaat met lengte ten hoogste  $f(|z|)$  en minimale editeerafstand van  $z[]$ . Kies  $f()$  minimaal – dus zo dat er in feite  $z[]$  en  $t[]$  bestaan waarin geen kortere beste benaderende matches zijn.*

Maar aan de andere kant zien jullie ook dat bij deze manier van doen veel deelresultaten altijd opnieuw berekend worden – de afstand voor dezelfde prefix wordt bv. voor elk woord berekend waarvan het een prefix is. En dan is natuurlijk duidelijk dat het veel sneller kan door opnieuw het principe van dynamisch programmeren toe te passen. Maar omdat er dan nog altijd  $O(m * n^2)$  editeerafstanden berekend en opgeslagen worden zal de complexiteit nog altijd ten minste  $O(m * n^2)$  zijn.

Wij gaan dus nog 1 stap verder: voor de  $j$ -de letter in de tekst (dus  $t[j-1]$ ) en een lengte  $i$  van een prefix slaan wij niet de afstanden van  $z[0], \dots, z[i-1]$  van alle deelwoorden van  $t[]$  op die op positie  $j-1$  eindigen, maar alleen de minimale editeerafstand van een deelwoord dat er eindigt.

Of precies:

$$D[i][j] = \min\{d((z[0], \dots, z[i-1]), (t[l], \dots, t[j-1])) \mid 0 \leq l \leq j\}$$

waarbij  $(z[0], \dots, z[-1])$  en  $(t[j], \dots, t[j-1])$  lege woorden zijn. Zie ook hier Oefening 77.

Ook hier zijn de eerste rij en de eerste kolom van  $D[][]$  gemakkelijk om te berekenen:

$D[i][0] = i$  voor alle  $i$  omdat alle  $i$  lettertekens verwijderd moeten worden.

$D[0][j] = 0$  voor alle  $j$  omdat de beste match van de lege string natuurlijk met de lege string is en de beste afstand dus 0 is.

**Volledig** analoog met het berekenen van de editeerafstand tussen twee strings kunnen wij nu voor  $i > 0$  en  $j > 0$  argumenteren:

$$D[i][j] = \min\{D[i-1][j] + 1, D[i][j-1] + 1, D[i-1][j-1] + g(i, j)\}$$

waarbij  $g(i, j) = 0$  als  $z[i-1] = t[j-1]$  en anders 1.

Als wij deze tabel ingevuld hebben, moeten wij alleen nog de kleinste waarde in de laatste rij zoeken (dus de kleinste afstand die het hele woord  $z[]$  heeft). Dat is de gezochte minimale afstand.

**Algoritme 14** (*Dynamisch programmeren om de kleinste editeerafstand van een deelwoord te berekenen*)

```
shortest_editeerafstand_deelwoord(z[], t[])
{
    // het wordt verondersteld dat z, t beide niet leeg zijn en dat de
    // lengte |z|, resp. |t| is. De kleinste editeerafstand van z met
    // een deelwoord van t wordt teruggegeven.

    for (i=0; i<=|z|; i++) D[i][0]=i;
    for (j=0; j<=|t|; j++) D[0][j]=0;

    for (i=1; i<=|z|; i++)
    {
        for (j=1; j<=|t|; j++)
        {
            // als k lettertekens beschouwd worden heeft het laatste
            // index k-1
            if (z[i-1]==t[j-1]) buffer=D[i-1][j-1];
            else buffer=D[i-1][j-1]+1;
            D[i][j] = min{buffer, D[i-1][j]+1, D[i][j-1]+1};
        }
    }
}
```

```

best=oneindig;
for (i=0;i<=|t|;i++) if (D[|z|][i]<best) best=D[|z|][i];
return best;
}

```

Merk op dat dit algoritme bijna identiek is aan Algoritme 13. Alleen de initialisatie van de eerste rij en eerste kolom verschillen en op het einde wordt nog het kleinste getal in de laatste rij gezocht. Wij hebben dus onmiddellijk:

**Lemma 9** *De kleinste editeerafstand tussen een string  $z[]$  van lengte  $n$  en een deelwoord van een string  $t[]$  van lengte  $m$  kan in tijd  $O(n * m)$  berekend worden.*

Maar ook hier is de implementatie en het besparen van geheugen belangrijk. Zie Oefening 89.

De werking van het algoritme voor het voorbeeld  $z[] = \text{examen}$  en  $t[] = \text{de\_armen\_vouwen}$  kan in de volgende tabel gezien worden.

	d e _ a r m e n _ v o u w e n														
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
e	1	1	0	1	1	1	1	0	1	1	1	1	1	1	0
x	2	2	1	1	2	2	2	1	1	2	2	2	2	2	1
a	3	3	2	2	1	2	3	2	2	2	3	3	3	3	2
m	4	4	3	3	2	2	2	3	3	3	3	4	4	4	3
e	5	5	4	4	3	3	3	2	3	4	4	4	5	5	4
n	6	6	5	5	4	4	4	3	2	3	4	5	5	6	5

**Oefening 84** *Pas Algoritme 14 toe op het woord  $z[] = \text{zoeken}$  en  $t[] = \text{koekjes\_kweken}$ . Bouw de tabel op en toon aan wat de beste match is en welke wijzigingen nodig zijn om **zoeken** in deze match te wijzigen.*

**Oefening 85** *Pas Algoritme 14 toe op het woord  $z[] = \text{test}$  en  $t[] = \text{beste\_tekst}$ . Bouw de tabel op en toon aan wat de beste matches zijn en welke wijzigingen nodig zijn om **test** in deze matches te wijzigen.*



**Oefening 86** Beschrijf een algoritme in pseudocode dat als  $z[]$ ,  $t[]$  en de tabel die door Algoritme 14 opgebouwd wordt ingevoerd worden de beste match teruggeeft.

**Oefening 87** Beschrijf een algoritme in pseudocode dat voor een woord  $z[]$  van lengte  $m$  en een tekst  $t[]$  van lengte  $n$  alle editeerafstanden van  $z[]$  met deelwoorden uit  $t[]$  berekent in tijd  $O(m * n^2)$ . Geef uitleg over de complexiteit van jouw algoritme.

**Oefening 88** Stel dat de kost van een verwijder-, toevoeg- en wisselbewerking om de ene string in de andere om te vormen niet constant 1 is maar dat elke van de 3 bewerkingen een individuele (positieve) kost heeft. Welke resultaten zijn nog altijd geldig en hoe kan je voor dit probleem een efficiënt algoritme schrijven?

**Oefening 89** Herschrijf de pseudocode van Algoritme 14. **Stel** daarbij dat het antwoord op de vraag in Oefening 77 ja is en let erop zo weinig mogelijk geheugen te gebruiken. Vooral hier kan de tekst natuurlijk veel te lang zijn om de hele matrix  $D[] []$  in het geheugen te houden! Welke van de twee for-lussen is beter de buitenste lus?

## 5 Compressiealgoritmen

Stel dat jullie de string

3,1415926535897932384626433832795028841971693993751058209749445923 willen comprimeren. Dan zouden jullie misschien zeggen dat dat een benadering van Pi is en dus voor de meeste toepassingen zeker 3,141592653 voldoende precies is (en gelijk hebben). Wat jullie dan doen, is compressie *met verlies*. Het resultaat van de compressie laat het niet toe de oorspronkelijke data te herstellen. Om te beslissen wat belangrijk is en wat niet moet je weten wat de data voorstelt en hoe hij geïnterpreteerd wordt. Als iemand bv. alleen in elk vijfde cijfer geïnteresseerd was, zou deze manier van compressie heel slecht zijn: je houdt overbodige informatie bij en gooit belangrijke informatie weg. Belangrijke formaten zoals jpg of mp3 zijn ook compressie met verlies: je kan de oorspronkelijke foto of het oorspronkelijke geluid niet reconstrueren. De formaten zijn erop gebaseerd dat geweten is hoe ons zicht en gehoor functioneren. Omdat het voor dergelijke formaten ook vooral belangrijk is te verstaan hoe de waarneming werkt (wat niet ons onderwerp is) zullen wij die niet bespreken maar ons alleen bezig houden met compressie zonder verlies – dus met compressiealgoritmen die het toelaten elke bit van het originele bestand te reconstrueren.

Een eerste vraag zou natuurlijk zijn of er een compressiealgoritme is dat elk bestand kan comprimeren en het antwoord is duidelijk *nee*: Omdat wij elk bestand willen reconstrueren moet een compressiealgoritme een bijjectie definiëren tussen de gecomprimeerde bestanden en de bestanden die gecomprimeerd worden. Als twee bestanden op hetzelfde bestand afgebeeld zouden worden, zou je het oorspronkelijke bestand niet kunnen reconstrueren! Maar hoe minder bytes hoe minder mogelijke bestanden – dus kan dat zeker niet. En even gemakkelijk kunnen wij zien dat als er ook maar één bestand is dat echt gecomprimeerd wordt (dus waar het resultaat korter is) dat er dan ook tenminste één bestand moet zijn dat langer wordt.

**Oefening 90** • *Bewijs expliciet dat als er een routine voor het comprimeren van bestanden is die een bestand echt comprimeert dat er dan ook ten minste één bestand is dat door deze routine langer gemaakt wordt.*

- *Een een beetje vage uitspraak: elk bestand kan gecomprimeerd worden. Schrijf deze uitspraak op een **precieze** manier met kwantoren (voor elk bestand bestaat...) en bewijs de precieze uitspraak of vind een tegenbeeld.*

Maar hoeveel bestanden van een gegeven lengte kan je met een gegeven compressiealgoritme comprimeren? Er zijn

$$\sum_{i=0}^n 256^i = \frac{1}{255} 255 \sum_{i=0}^n 256^i = \frac{1}{255} (256^{n+1} - 1)$$

bestanden met ten hoogste  $n$  bytes. Als wij nu bestanden met  $n + 1$  bytes willen comprimeren (waarbij wij altijd ten minste één byte willen winnen) dan is het beeld van een gecomprimeerd bestand een bestand met ten hoogste  $n$  bytes. Omdat er  $256^{n+1}$  bestanden met  $n + 1$  bytes maar alleen maar ongeveer  $256^{n+1}/255$  bestanden met ten hoogste  $n$  bytes zijn, kan dus **ten hoogste**  $1/255$ de van de bestanden gecomprimeerd worden – en daarvoor moeten wij dan nog *betalen* omdat andere bestanden langer worden! Dat klinkt niet goed. . .

Maar onze ervaring is helemaal anders: de meeste bestanden waarop wij compressiealgoritmen zoals gzip, zip, bzip2, etc. toepassen, kunnen gecomprimeerd worden – in sommige gevallen zelfs heel sterk. De reden lijkt op de reden waarom metaheuristieken werken: als wij het maximum van een toevallige functie zouden moeten vinden, zouden metaheuristieken nutteloos zijn. Zij werken goed omdat de problemen waarmee mensen meestal bezig zijn structuur vertonen – ze zijn helemaal niet toevallig! En hier is het hetzelfde: de bestanden die wij normaal willen comprimeren zijn niet toevallig. Het zijn bv. tekstbestanden met woorden uit een zekere taal. In dergelijke bestanden komen sommige letters veel vaker voor dan anderen – en dat is een structuur die toevallige bestanden niet tonen. Dergelijke structuren maken het mogelijk bestanden te comprimeren. Alle technieken die wij zullen zien, zullen slecht presteren op bestanden die gewoon toevallige bits bevatten. En wij hebben net getoond dat dat ook niet anders kan!

**Oefening 91** *Wij hebben aangetoond dat je met een gegeven compressiealgoritme ten hoogste ongeveer  $1/255$  van alle bestanden kan comprimeren. Maar dan hebben wij aanvaard dat het misschien een compressie van maar 1 byte is.*

- *Bereken hoe groot het aandeel is van bestanden die door een vast compressiealgoritme met ten minste 100 bytes gecomprimeerd kunnen worden.*
- *Geef een bovengrens voor de kans dat een vast compressiealgoritme een toevallig bestand van 5MB of meer om ten minste de helft kan comprimeren*

**Oefening 92** *Wij hebben net gezien dat je maar 1/255ste van alle bestanden kan comprimeren. Stel nu dat elke computer ongeloofelijk veel ruimte ter beschikking heeft en op voorhand al grote delen kan opslaan die in bestanden kunnen voorkomen (dus heel grote voorgëimplementeerde woordenboeken – bv. alle bestanden t.e.m. 500 bytes) zodat je op de volgende manier kan comprimeren: je slaat gewoon de index van de onderdelen op en stuurt die door. De index van de woorden kan je natuurlijk nog eens met andere compressiealgoritmen bewerken. Bepaal voor deze manier van comprimeren een bovengrens voor het aantal bestanden dat je kan comprimeren. Resultaten en bewijzen uit de les mogen geciteerd worden en moeten niet herhaald worden.*

Een eerste aanzet:

Stel dat je niet meer een vaste lengte voor alle codewoorden (bv. bytes) wil hebben, maar variabele lengte van de codewoorden toelaat. Je kan bv. bytes coderen door nullen in het begin weg te laten. Dus bv. 00011111 als 11111. Dat lijkt al ruimte te besparen, maar het probleem is natuurlijk dat je op deze manier niet meer kan herkennen waar het ene woord stopt en het volgende begint als je ze in een bestand achter elkaar schrijft. Waar zijn bv. de grenzen in 1111101101010111010111?

Eén manier om dit op te lossen is het gebruik van prefix-codes:

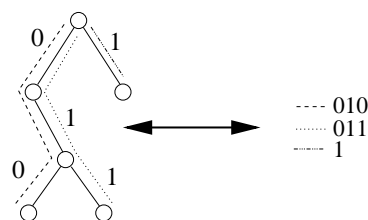
**Definitie 7** *Een verzameling  $M$  van woorden (met tekens uit een zeker alfabet) heet een prefix-code als voor elk woord  $w = w_1, w_2, \dots, w_m$  geen woord  $w' = w'_1, w'_2, \dots, w'_n$  bestaat met  $n > m$  en  $w_i = w'_i$  voor alle  $1 \leq i \leq m$ .*

Je kan dus gemakkelijk herkennen wanneer je het einde van een woord hebt bereikt – gewoon omdat als je doorgaat er geen element in de verzameling zit dat door de langere code gecodeerd zou kunnen zijn! In de toekomst zal het alfabet waarover wij het hebben gewoon  $\{0, 1\}$  zijn – wij hebben het dus over woorden van bits.

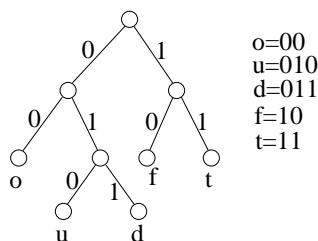
Prefixcodes kan je altijd in de vorm van een boom voorstellen waar de bogen gelabeld zijn met letters uit het alfabet. De codewoorden corresponderen dan met de labellings langs de paden van de wortel naar de bladeren.

De prefixeigenschap correspondeert dan met het bereiken van een blad: er is geen langer woord met hetzelfde begin omdat het blad anders een kind had (en dan zou het geen blad zijn).

De rijen van nullen en enen staan normaal voor tekens – zoals bv. de ASCII tekens beschreven worden door rijen van nullen en enen met lengte 8. In een boom kunnen wij dat aanduiden door de letters aan de bladeren te schrijven. Wat zou bv. de string 00010011 coderen als wij de code nemen die in Figuur 17 beschreven is?



Figuur 16: Een kleine prefixboom en de bijbehorende codewoorden.



Figuur 17: Een kleine prefixboom, de bijbehorende codewoorden en tekens.

Je begint met het eerste bit en volgt de takken in de boom. Als je een blad bereikt, schrijf je het teken dat bij dit blad behoort op. Dan is één codewoord gedecodeerd en beginnen wij met de volgende bit opnieuw bij de wortel om het volgende teken op dezelfde manier te decoderen. Als je dat toepast tot alle bits gelezen zijn, heb je het hele woord gedecodeerd.

In het voorbeeld lezen wij eerst 00 en hebben het blad *o* bereikt. Dan beginnen wij opnieuw met de wortel en bereiken het blad *u* na 010, beginnen opnieuw en bereiken *d* na 011. De string die er gecodeerd was was dus het woord *oud*.

**Oefening 93** *Wat is het juiste antwoord op de vraag wat de string 100001011 codeert als wij de code nemen die in Figuur 17 beschreven is?*

Als wij deze manier van coderen toepassen, kunnen de codes voor verschillende tekens ook een verschillende lengte hebben. Als wij een code willen hebben die zo kort mogelijk is, is het dus verstandig korte codes voor tekens te gebruiken die vaak opduiken en de langere codes voor tekens die minder vaak gebruikt worden. Om een voordeel aan deze techniek te hebben, moet je wel bestanden hebben waar de tekens niet *toevallig met gelijke kans* verdeeld zijn, maar die structuur vertonen – waar sommige tekens dus veel vaker opduiken dan anderen.

Als wij een alfabet  $A$  hebben en voor  $x \in A$  het aantal keren dat  $x$  in het bestand zit dat wij willen comprimeren  $a(x)$  is en de lengte van zijn code  $l(x)$  dan is de lengte van het gecodeerde bestand  $\sum_{x \in A} (a(x) * l(x))$  bits. Wij

hebben geen invloed op de functie  $a()$  – het bestand is gegeven – maar wij kunnen de codering kiezen en dus invloed op de  $l(x)$  uitoefenen. Ons doel is dus de codering op een manier te kiezen dat  $\sum_{x \in A} (a(x) * l(x))$  minimaal is.

## 5.1 Huffman codering

(Vanaf 2008 niet meer deel van de les.)

Huffman codering is een gretige manier om een code te construeren. Maar terwijl gretige algoritmen vaak alleen maar benaderende resultaten opleveren kunnen wij in dit geval aantonen dat het resultaat optimaal is.

Huffman codering werkt op de volgende manier:

- Doorloop eerst het bestand en bepaal voor elk teken  $x$  de frequentie  $a(x)$ .

Nu zullen wij de coderingsboom stap voor stap opbouwen. De boom zal niet uniek bepaald zijn, maar alle bomen die je op deze manier krijgt, zijn even goed:

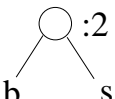
### Algoritme 15 Huffman codering:

- *maak een verzameling van bomen aan waarin elke letter met een strikt positieve frequentie een boom met één top vormt.*
- *herhaal de volgende stap tot er nog maar één enkele boom in de verzameling zit:*
  - *Voeg de twee bomen  $X, Y$  met de kleinste frequentie samen tot één boom  $Z$  door van de wortels van  $X$  en  $Y$  de twee kinderen van één nieuwe wortel te maken. De frequentie  $a(Z)$  wordt gedefinieerd als  $a(Z) = a(X) + a(Y)$ .*

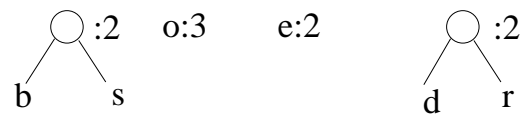
*De lettertekens zijn achteraf de bladeren van deze boom.*

Als voorbeeld zullen wij nu de korte tekst **boosdoer** coderen waar de lesgever een  $n$  is vergeten:

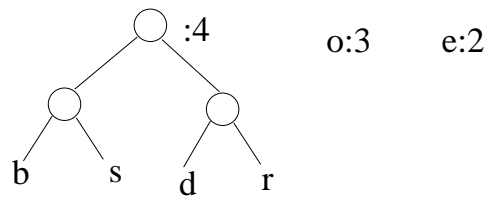
b:1    o:3    e:2    s:1    d:1    r:1

eerste stap:  o:3    e:2    d:1    r:1

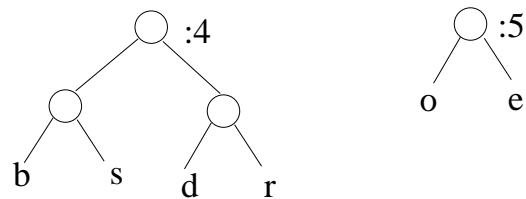
tweede stap:



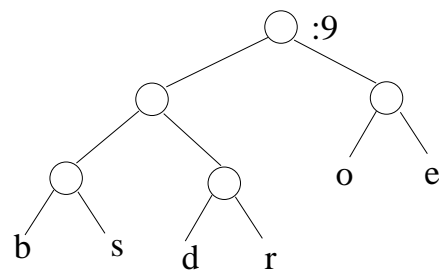
derde stap:  
(niet uniek!)



vierde stap:



vijfde stap:



De codes zijn dus:

**b:** 000

**o:** 10

**e:** 11

**s:** 001

**d:** 010

**r:** 011

**Oefening 94** Bereken de Huffman codering voor bananenboom.

Inderdaad is dat een gretig algoritme, en wij hebben gezien dat gretige algoritmen soms alleen benaderingen berekenen – en in sommige gevallen zelfs slechte benaderingen. Maar in dit geval krijgen wij inderdaad de optimale oplossing:

**Stelling 10** *Voor een gegeven alfabet  $x_1, \dots, x_k$  en frequenties  $a(x_i)$  berekent het Huffman algoritme een prefixcode (boom)  $T$  waarvoor*

$$Q(T) = \sum_{i=1}^k l(x_i) * a(x_i)$$

*minimaal is, waarbij  $l(x_i)$  de lengte van het codewoord voor  $x_i$  is.*

**Bewijs:** Let op: voor het Huffman algoritme zijn de letters in het alfabet natuurlijk kleine boompjes. Wij zullen het dus altijd over bomen hebben.

**Opmerking:** In een optimale boom met ten minste 3 toppen heeft elke top ofwel 0 ofwel 2 kinderen. Toppen met maar 1 kind zou je gewoon kunnen verwijderen (en de ouder en het kind met elkaar verbinden of de top alleen maar verwijderen als het de wortel was) en het resultaat zou een betere boom zijn: voor alle bladeren  $x$  blijft  $l(x)$  gelijk of wordt kleiner. Dat betekent ook dat elke top (behalve de wortel) een broer heeft (een andere top met dezelfde ouder).

Wij gebruiken inductie in het aantal stappen waarin bomen samengevoegd worden (dat zijn er altijd  $k - 1$  als  $k$  het aantal boompjes in het begin van het algoritme is). Als er ten hoogste één stap is, is het resultaat uniek (er is maar één mogelijke optimale boom) – dat is dus triviaal.

Stel nu dat bewezen is dat het Huffman algoritme voor  $s - 1$  stappen altijd een optimale boom bouwt en dat nu  $s$  stappen nodig zijn. Er zijn dus  $s + 1$  letters. Stel bovendien dat het algoritme als de boom  $T_{s+1}$  voor  $s + 1$  lettertekens opgebouwd wordt eerst de laatste 2 letters in een boompje  $T_{x_s, x_{s+1}}$  samenvoegt (anders moeten de letters gewoon hernoemd worden) en er de frequentie  $a(T_{x_s, x_{s+1}}) = a(x_s) + a(x_{s+1})$  aan toekent. Deze twee hebben dus de kleinste frequenties (maar ze zijn misschien niet de enigen met deze eigenschap).

Dan bouwt het algoritme een optimale boom  $T_s$  voor de verzameling van bomen  $x_1, \dots, x_{s-1}, T_{x_s, x_{s+1}}$  als wij  $T_{x_s, x_{s+1}}$  als één letterteken beschouwen. Stel nu dat de boom  $T_{s+1}$  die voor  $x_1, \dots, x_{s+1}$  gebouwd wordt – dat is dezelfde boom maar met  $T_{x_s, x_{s+1}}$  als boom van diepte 1



in plaats van een enkele top – niet optimaal is maar dat er een betere optimale boom  $T_{opt}$  bestaat.

**Opmerking:** Wij mogen stellen dat  $x_s, x_{s+1}$  broers zijn in  $T_{opt}$  en op het laagste niveau zitten:

Neem twee broers  $y, y'$  op het laagste niveau. Stel dat  $y \notin \{x_s, x_{s+1}\}$  en  $x \in \{x_s, x_{s+1}\}, x \notin \{y, y'\}$ . Dan kunnen wij  $x$  met  $y$  wisselen en krijgen en ten minste even goede boom: Wij hebben  $l(y) \geq l(x)$  (omdat  $y$  op het laagste niveau zit) en  $a(y) \geq a(x)$  (omdat  $x$  één van de twee toppen met laagste frequentie is).

Voor de nieuwe boom  $T'_{opt}$  na het wisselen zou dus gelden:

$$Q(T'_{opt}) = Q(T_{opt}) - ((a(x)l(x)) + a(y)l(y)) + ((a(x)l(y)) + a(y)l(x)) = Q(T_{opt}) - (a(y) - a(x))(l(y) - l(x)) \leq Q(T_{opt})$$

De nieuwe boom zou dus ten minste even goed zijn. Wij kunnen dus  $x_s, x_{s+1}$  met andere toppen wisselen zodat aan de voorwaarden uit de opmerking voldaan is.

Als wij in deze boom nu  $x_s, x_{s+1}$  en de ouder door één top  $z$  met frequentie  $a(x_s) + a(x_{s+1})$  vervangen dan geldt voor de zo ontstane boom  $T_{opt,s}$ :

$$Q(T_{opt,s}) = Q(T_{opt}) - (a(x_s) + a(x_{s+1})) < Q(T_{s+1}) - (a(x_s) + a(x_{s+1})) = Q(T_s)$$

Maar dan zou  $T_s$  niet optimaal zijn wat een tegenstrijdigheid met de inductiehypothese is.

■

Je kan een tekst dus coderen door eerst de boom op te slaan die de code beschrijft en daardoor vastlegt hoe je de tekst moet decoderen, en achteraf de code volgens deze boom. Over de manier waarop je de boom het best opslaat, zullen wij het hier niet hebben – compressie is alleen maar nuttig voor grote bestanden en dan is het kleine en voor een gegeven alfabet constante deel voor de boom zonder veel belang.

Wat wij Huffman codering genoemd hebben kan je ook *statische Huffman codering* noemen. Wij gaan ervan uit dat je de frequenties van elk letterteken in het bestand kent. Aan de ene kant betekent dat natuurlijk dat je het bestand twee keer moet lezen – één keer om de frequenties te bepalen en één keer om het te comprimeren. Dat is natuurlijk niet ideaal. Maar soms is het zelfs onmogelijk: als je geen bestand maar een stream van data wil comprimeren dan moet dat *online* gebeuren – je **kan** de stream niet helemaal lezen

en dan nog een keer opvragen. Daarvoor kan je adaptieve Huffman codering gebruiken. Het verschil is gewoon dat je een Huffman boom online opbouwt en wijzigt als de frequenties veranderen (het lijkt een beetje op het updaten van zoekbomen). Aan de kant waar je de code moet decoderen wordt op dezelfde manier dezelfde boom opgebouwd zodat je de code kan decoderen. Zo moet de boom die de code beschrijft ook niet opgeslagen worden. Details van adaptieve Huffman codering zullen wij hier niet bespreken – die kunnen ook gemakkelijk zelf uitgewerkt worden – maar in plaats daarvan nog andere compressiealgoritmen die nog beter presteren en ook op streams toepasbaar zijn.

**Oefening 95** *Construeer de Huffman code van de volgende tekst:*

`zaken_nemen_geen_keer`

**Oefening 96** *Stel een manier voor om statische Huffman codering voor streams te gebruiken. Wat zijn de voordelen en nadelen in vergelijking met statische Huffman codering als je die op de hele verzameling van data kan toepassen?*

**Oefening 97** *Een geautomatiseerd systeem om goederen in een opslagplaats op te slaan kan ook onderafdelingen van het magazijn samenvoegen – maar nooit meer dan 2 tegelijk.*

*De tijd die het nodig heeft om deze onderafdelingen samen te vatten is  $c * (d_1 + d_2)$  als  $d_1$  het aantal goederen in afdeling 1 is en  $d_2$  het aantal goederen in afdeling 2.*

*Een groot bedrijf wil nu al zijn  $n$  onderafdelingen tot één grote afdeling samenvatten. De hoeveelheden  $d_1, \dots, d_n$  van goederen in de afdelingen zijn gekend.*

*Geef een  $O(n \log n)$  algoritme dat een volgorde bepaalt waarop de onderafdelingen zo snel mogelijk samengelegd kunnen worden.*

**Oefening 98** *Gegeven een vast compressiealgoritme  $A$ . De taak is een bestand van  $n$  MByte te vinden dat door het algoritme **niet** gecomprimeerd kan worden.*

*Geef een algoritme dat zo'n bestand construeert. Welk soort algoritme zou je voorstellen? Stel dat het algoritme  $A$  in tijd  $O(n)$  werkt. Hoeveel tijd zou je verwachten dat jouw algoritme vraagt – kan het even snel?*

**Oefening 99** *Bediscussieer de manier waarop je tijdens de Huffman codering altijd de twee boompjes met de minimale frequentie vindt. Welke data-structuren zou je kunnen toepassen?*

*Stel dat je op de volgende manier werkt:*

*Je houdt twee lijsten bij: lijst (a) bevat in het begin de tekens samen met hun frequenties – dat zijn dus ook kleine boompjes. Je moet de lijst in stijgende volgorde van de frequentie sorteren. De lijst (b) gaat de gebouwde boompjes bevatten en is in het begin leeg.*

*Als je nu moet beslissen welke twee boompjes je moet samenvoegen om de volgende boom te bouwen, kies je 2 keer het kleinste van de elementen die in de lijsten vooraan staan en verwijder je het uit zijn lijst. Dan vorm je daarvan de nieuwe boom en schrijf je hem op het einde van lijst (b).*

*Toon aan dat dit algoritme correct werkt – dus: dat je op deze manier altijd de 2 bomen met de kleinste frequenties samenvoegt.*

*Hoeveel tijd vraagt dit algoritme voor  $n$  tekens met gegeven frequenties ?*

Wij hebben *in zekere zin* bewezen dat *Huffman codering optimaal is*. Maar wat betekent dat precies? Is er geen betere prefixcode om iets op te slaan? Daarvoor moeten wij precies kijken wat de voorwaarden van de stelling waren: er is geen betere prefixcode voor een **gegeven** alfabet. Maar stel bv. dat in de tekst het letterteken “o” altijd gevolgd wordt door “e”. Dan zou het toch verstandig zijn niet “o” en “e” voor de code te gebruiken maar “oe”! En zelfs in sommige gevallen waar er “o”, “e” en “oe” in de code zitten, zou het kunnen dat het nuttig is alle drie te coderen.

En dat geldt natuurlijk ook voor langere woorden. Je kan dus het alfabet uitbreiden door kleine deelstrings als lettertekens te beschouwen. In het extreme geval zou je natuurlijk de hele tekst als één letterteken kunnen beschouwen en de code is alleen maar “1”. Maar dan is de benadering dat de ruimte om de code op te slaan verwaarloosbaar is zeker niet meer van toepassing...

De basisideeën om betere compressiealgoritmen te ontwikkelen dan Huffman codering zijn

- Gebruik een groter alfabet dan alleen maar lettertekens – gebruik ook strings!
- Gebruik een woordenboek dat dynamisch kan veranderen.

**Oefening 100** *Stel dat je niet één groot bestand moet sorteren, maar  $n$  al bestaande gesorteerde bestanden  $b_1, \dots, b_n$  tot één groot gesorteerd bestand moet samenvoegen. Voor  $1 \leq i \leq n$  bevat bestand  $b_i$   $l_i$  sleutels en je kan de bestanden niet tegelijk mergen maar je kan altijd alleen maar 2 bestanden mergen tot een nieuw bestand. Als de oude bestanden  $l_i$  en  $l_j$  sleutels bevatten, bevat het nieuwe bestand dus  $l_i + l_j$  sleutels en de kost van deze mergebewerking is  $l_i + l_j$ .*

- *Geef een voorbeeld dat toont dat de kost om alle bestanden te mergen afhankelijk is van de volgorde waarop je de bestanden mergt.*
- *Beschrijf een efficiënt algoritme om de optimale (goedkoopste) volgorde te bepalen.*
- *Wat is de complexiteit van dit algoritme?*
- *Toon aan dat jouw algoritme inderdaad de optimale volgorde vindt.*

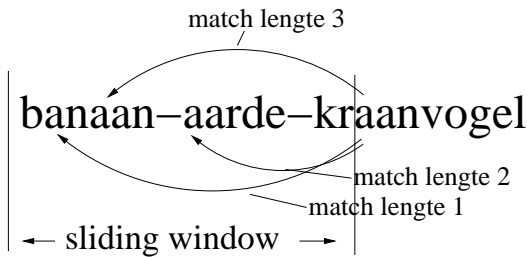
## 5.2 LZ77

(Abraham Lempel en Jacob Ziv 1977)

Het LZ77 algoritme is de basis van sommige algoritmen die in de praktijk heel goed presteren en vaak gebruikt worden – zoals bv. zip en gzip.

**Het idee:** Het idee is dat je de tekst zelf als jouw woordenboek gebruikt. Een woordenboek is natuurlijk alleen maar nuttig als het ook bij de context past (bv. zou een Nederlands woordenboek nauwelijks nuttig zijn voor een Engelse tekst), maar de tekst als woordenboek past natuurlijk ideaal bij zijn eigen context! Dus: aan de ene kant heb je de tekst toch al nodig en aan de andere kant heb je dan ook automatisch de woorden die je nodig hebt. Natuurlijk kan je niet de hele tekst bijhouden (als het bv. een stream is) – dus gebruik je altijd alleen maar een eindig deel van de tekst – bv. de laatste 32 KB die je gezien hebt. Dit eindige deel noemen wij het *sliding window*. Dit sliding window is in het begin leeg en bevat tijdens het algoritme altijd de laatste  $G$  bytes waarbij  $G$  een getal is dat je op voorhand vastlegt. Zo heb je wel minder woorden ter beschikking maar als de tekst qua structuur verandert heb je de actuele woorden toch al ter beschikking! In de code houd je dan altijd bij waar het gecodeerde woord in het sliding window begint en hoe lang het is. Wat *het woord* is waarnaar je verwijst hangt ervan af wat de langste match in het sliding window is. Maar je moet er natuurlijk ook voor zorgen dat in het geval dat er helemaal geen match in het sliding window is het volgende letterteken gecodeerd kan worden.

Het algoritme werkt als volgt:



neem langste match: positie 3 lengte 3

### Algoritme 16 LZ77:

Stel dat het bestand dat gecodeerd moet worden de rij van lettertekens  $x_0, x_1 \dots$  is. Als wij het over een positie in het bestand hebben dan bedoelen wij de index. Analooog gebruiken wij positie voor de lettertekens die op dat moment in het sliding window staan – daar beginnen wij met het eerste teken in het sliding window (om het even wat de index van dit teken in het bestand is) en geven die het nummer 0.

- Houd altijd een sliding window van de laatste  $G$  bytes bij waarbij  $G$  een op voorhand vastgelegd getal is.
- Als de tekst vanaf  $x_k$  gecodeerd moet worden dan wordt de langste match gezocht voor een woord dat op positie  $k$  begint. De match moet met een letterteken in het sliding window beginnen – maar het einde van de match mag er wel buiten liggen.

Stel dat de index in het sliding window van het begin  $p$  is, de lengte  $l$  en het letterteken dat op de langste match volgt  $x$ . Als er geen match is, neem gewoon  $p = 0$  maar omdat  $l = 0$  doet het er niet echt toe – ook andere waarden voor  $p$  zouden dan werken.

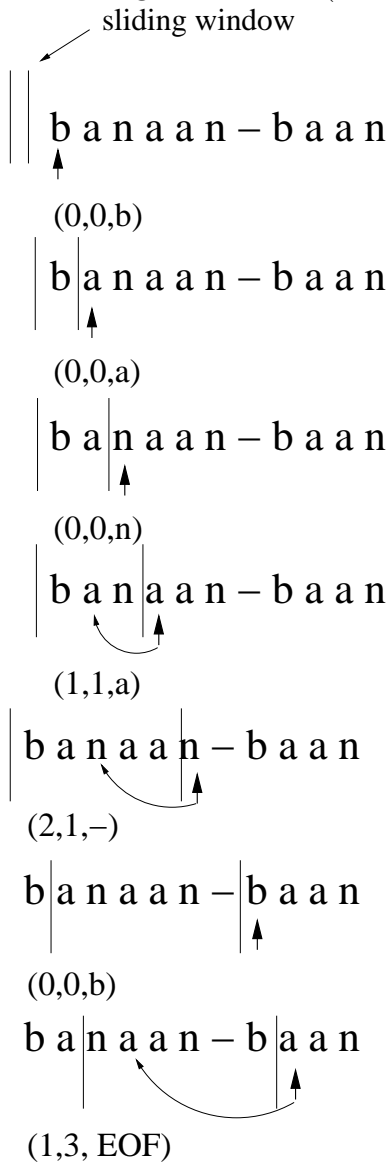
Dan wordt het 3-tal  $(p, l, x)$  als code geschreven.

- Het volgende teken vanaf waar gecodeerd moet worden, is op positie  $k + l + 1$  in het bestand.

Daarbij kan je natuurlijk de positie op verschillende manieren beschrijven – bv. door het aantal tekens vanaf het begin van het sliding window of vanaf het einde van het sliding window, etc.

Je kan ook getallen en lettertekens als aparte lettertekens voor een alfabet beschouwen en dan ofwel de positie en de lengte opslaan ofwel het niet gevonden teken (dus  $(p, l)$  ofwel  $x$ ). Omdat het verschil tussen getallen en lettertekens herkend kan worden is dat geen probleem. Of je kan alleen maar het letterteken schrijven als de lengte 0 is ofwel...

Al deze manieren van doen zijn mogelijk en verwezenlijken het idee van LZ77. Wij zullen nu een voorbeeld zien waar wij de codering met 3-tallen gebruiken. In de les gebruiken wij alleen deze codering waarbij we de positie vanaf de linkerkant tellen en het eerste teken de positie 0 krijgt. Stel dat de *tekst* *banaan-baan* gecodeerd moet worden en dat de grootte van het sliding window 6 is (om iets te tonen moet het natuurlijk zo klein zijn. . .)



Deze output kan dan nog Huffman gecodeerd worden zodat er bv. rekening mee gehouden kan worden dat sommige lengten vaker voorkomen dan andere. Daarbij heb je verschillende mogelijkheden: je kan de output als een reeks van bytes zien – dus ook de lengten gewoon als strings van cijfers. Maar de

symbolen in jouw Huffman-boom zouden ook de hele drietallen kunnen zijn. Je kan ook niet alleen één Huffman boom gebruiken maar bv. twee – één voor de getallen die posities en lengten beschrijven en één voor de lettertekens – of zelfs drie – één voor de posities, één voor de lengten en één voor de lettertekens. Omdat je elke keer als je een nieuw symbool wil lezen, weet of het een lengte, een positie of een letterteken moet zijn, weet je welke boom je moet gebruiken. In het geval van de interpretatie van de output als gewone tekst van bytes heeft dat zeker voordelen. Je kan dan bv. dezelfde korte bitstring voor verschillende tekens gebruiken in de verschillende omstandigheden. In het geval dat je de drietallen als tekens beschouwt, heb je dan wel kortere codes voor elk teken – maar er zijn ook twee of drie codes die gelezen moeten worden terwijl anders een heel drietal door één code beschreven wordt. Het is dus niet op voorhand duidelijk wat in dit geval beter is. Een voordeel zou misschien kunnen zijn dat meer rekening wordt gehouden met de individuele hoeveelheden van de verschillende symbolen in het drietal. Terwijl zekere posities of lengtes misschien vaker voorkomen dan andere, is dat misschien minder het geval als zich hele drietallen moeten herhalen. Of dat echt zo is en voldoende winst oplevert om meer codes te rechtvaardigen is een goede programmeeroefening! De programma's *DEFLATE* en *gzip* gebruiken bv. twee bomen: één voor de posities van de matches en één voor de rest (lengtes en het volgende teken) (waarbij nog verschillende optimalisaties aan bod komen) – dat is een mogelijke oplossing die goede resultaten oplevert. Tijdens het decoderen wordt hetzelfde sliding window opgebouwd – altijd als een code naar het woordenboek verwijst staat het gezochte woord dus ook ter beschikking. Ook dit kan het gemakkelijkst met een voorbeeld duidelijk gemaakt worden. Wij stellen ook hier dat de lengte van het sliding window maar 6 is:

(0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 ↙ sliding window  
 | a |  
 (0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 | aan |  
 (0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 | aanv |  
 (0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 a | anvaar |  
 (0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 aa | nvaard |  
 (0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 aan | vaardb |  
 (0,0,a) (0,1,n) (0,0,v) (0,2,r) (0,0,d) (0,0,b) (1,3,EOF)  
 ↑  
 aanvaa | rdbaar |

Natuurlijk zijn deze voorbeelden weinig overtuigend – de codes zijn normaal veel langer dan het door ons gecodeerde woord en ook het sliding window maakt natuurlijk weinig kans overeenkomsten te vinden als het zo klein is! Maar realistische voorbeelden kan je natuurlijk niet met de hand uitwerken. . .

### Oefening 101 *Codeer de teksten*

- `swiss_miss_missing`
- `een_benen_been`
- `ananas-kanaal`

met LZ77 (zonder het resultaat achteraf nog met het Huffman algoritme te coderen). Gebruik een sliding window van grootte 8. Denk eraan dat matches in de sliding window moeten beginnen maar niet noodzakelijk eindigen!



**Oefening 102** *In dit voorbeeld hebben wij de code achteraf niet nog Huffman gecodeerd. Decodeer de volgende twee codes met een sliding window met grootte 6:*

- $(0, 0, e)(0, 1, n)(0, 0, d)(1, 2, b)(0, 4, EOF)$
- *Als je het volgende woord decodeert gebeurt er iets raars – maar dat kan:*  
 $(0, 0, b)(0, 0, l)(0, 0, a)(0, 9, EOF)$

**Oefening 103** *Welk algoritme voor string matching zou je gebruiken om een langste match te vinden. Geef redenen voor jouw keuze.*

**Oefening 104**

- *Comprimeer de volgende tekst met het LZ77 algoritme en pas achteraf Huffman met twee bomen toe. Gebruik een sliding window dat duidelijk groter is dan deze kleine tekst.*  
`abracadabra_simsalabim`
- *decodeer achteraf jouw code om te zien of je inderdaad hetzelfde resultaat krijgt.*

In de praktijk presteren op LZ77 gebaseerde algoritmen heel goed qua compressiefactor. Maar een probleem is de tijdscomplexiteit als je met grote sliding windows werkt. Natuurlijk is het zo dat hoe groter het sliding window is, hoe meer tijd het vinden van de langste match vraagt. Een sliding window heeft geen structuur – zoals een boom – die het je gemakkelijker maakt om te zoeken. Een dergelijke structuur opbouwen is ook niet zo eenvoudig omdat de sliding window altijd verandert...

Het decoderen is altijd efficiënt – om het even hoe groot het sliding window is. Als je de posities op die manier van de voorbeelden beschrijft, moet het sliding window tijdens het decoderen even groot zijn als tijdens het coderen. Als je de posities vanaf de rechterkant van het sliding window gebruikt, moet het niet even groot zijn maar ten minste even groot.

### 5.3 LZW

(Abraham Lempel, Jacob Ziv en Terry Welch 1984)

LZW is gebaseerd op het LZ78 algoritme dat Lempel en Ziv in 1978 voorstelden. Het grote verschil met LZ77 (en de reden waarom wij het hier voorstellen – wij willen tenslotte vooral nieuwe ideeën zien) is dat hier een

echt woordenboek wordt gebruikt en niet – zoals in LZ77 – de tekst zelf als woordenboek.

Natuurlijk is het woordenboek ook hier gebaseerd op de actuele tekst om alleen maar *nuttige* woorden in het woordenboek te hebben. Hoe groter het aantal woorden hoe langer de codes om de woorden te beschrijven – het is dus belangrijk zo weinig mogelijk overbodige woorden in het boek te hebben. De klemtoon ligt hier ook op de mogelijkheid efficiënt in het woordenboek te kunnen zoeken.

**Coderen:** Eerst zullen wij het principe uitleggen voordat wij de details beschrijven:

Je wil teksten van bytes coderen in codewoorden van een op voorhand vastgelegde lengte – meestal 12 bits. Je kan dus getallen  $0 \dots 4095$  voorstellen. Het woordenboek zal maximaal 4096 woorden bevatten die door deze nummers worden voorgesteld. In het begin vul je de *gewone bytes* als woorden in. Zij krijgen als nummer het getal dat ze voorstellen – dat is dus  $0, \dots, 255$ . De rest van het woordenboek zal nu dynamisch opgebouwd worden. Daarbij voeg je elke keer dat je een nieuw woord tegenkomt – dat is dan een woord dat 1 letterteken langer is dan de tot nu toe langste deelstring van dit woord in het woordenboek – dit woord toe aan jouw woordenboek. Dat blijf je doen totdat het woordenboek volzit. Je gebruikt de code van een woord dus pas de tweede keer dat je het woord tegenkomt – de eerste keer gebruik je het woord min het laatste teken en voeg je het woord toe aan het woordenboek. Maar natuurlijk is dat alleen maar een informele beschrijving. De details zie je in de volgende pseudocode:

voeg eerst de lettertekens toe aan jouw woordenboek  
en geef de nummers vanaf 0 in volgorde

```
//w is ons woord
w= de lege string

while (w!=EOF)
{
    k= volgende teken
    if (wk al in woordenboek)
        w=wk
    else
        { output code van w
          voeg wk toe aan het woordenboek
          geef wk de volgende code
```

```

        w=','k','
    }
}

```

Een voorbeeld toont hoe het algoritme werkt:

```

b a n a a n a a n v a l
↑      in boek

b a n a a n a a n v a l
↑ niet in boek   output code(b)=98
  voeg ba toe aan woordenboek met nummer 256

b a n a a n a a n v a l
↑ niet in boek   output code(a)=97
  voeg an toe aan woordenboek met nummer 257

b a n a a n a a n v a l
↑ niet in boek   output code(n)=110
  voeg na toe aan woordenboek met nummer 258

b a n a a n a a n v a l
↑ niet in boek   output code(a)=97
  voeg aa toe aan woordenboek met nummer 259

b a n a a n a a n v a l
in boek  ↑

b a n a a n a a n v a l
niet in boek  ↑   output code(an)=257
  voeg ana toe aan woordenboek met nummer 260

```

**b a n a a n a a n v a l**

in boek    ↑

**b a n a a n a a n v a l**

niet in boek    ↑    output code(aa)=259

voeg aan toe aan woordenboek met nummer 261

**b a n a a n a a n v a l**

niet in boek    ↑    etc.....

Als het woordenboek volzit kan je verschillende dingen doen:

- begin gewoon opnieuw met een nieuw woordenboek
- ga door met dit woordenboek tot het inefficiënt blijkt te zijn (weinig lange matches)
- verwijder het woord met meer dan 1 letterteken uit het woordenboek dat je het langst niet meer gebruikt hebt (als code of als prefix van een code).

Al deze manieren van werken worden toegepast (GIF, unix compress, British Telecom Standard) en jullie hebben zeker ook zelf interessante ideeën wat je zou kunnen doen...

### **Decoderen:**

Tijdens het decoderen wordt hetzelfde woordenboek opgebouwd als tijdens het coderen. Ook hier worden eerst de lettertekens aan het woordenboek toegevoegd. Elke keer dat een code gelezen wordt, zit die al in het woordenboek. Om dat te garanderen wordt tijdens het coderen een woord pas door zijn code beschreven als wij het woord voor de tweede keer tegenkomen. Dan wordt dit woord samen met het eerste letterteken van de volgende code aan het woordenboek toegevoegd. En hier hebben wij een probleem: Het volgende woord kan inderdaad de tweede keer zijn dat wij het woord tegenkomen dat wij net aan het woordenboek willen toevoegen.

Voorbeeld: Codeer **aaaaaaaa...**

De code begint met 97 voor **a** en **aa** wordt aan het woordenboek toegevoegd (code 256). Dan wordt 256 als code geschreven en **aaa** toegevoegd (code 257) etc.

Als wij nu decoderen, schrijven wij eerst **a** en dan willen wij **a** gevolgd door het eerste letterteken van het woord met code 256 aan het woordenboek toevoegen – maar nummer 256 willen wij net op dit moment bepalen...

Er is dus een probleem als (en slechts als) een woord onmiddellijk nadat het aan het woordenboek wordt toegevoegd als code wordt geschreven. Dan wordt het tijdens het decoderen niet teruggevonden. Maar wij weten al hoe het woord begint: omdat in dit geval het vorige woord alle lettertekens behalve het laatste bevat, begint het nieuwe woord met hetzelfde teken. En wij hebben alleen het eerste teken nodig om het nieuwe woord te bepalen dat toegevoegd moet worden.

Ook hier zie je de details het best in pseudocode waarbij **[woord,teken]** voor de string staat die met de string **woord** begint en daarop het teken **teken** volgt.

```
voeg eerst de lettertekens toe aan jouw woordenboek
en geef de nummers vanaf 0 in volgorde

// eerst de eerste code lezen

c=gelezen code // een getal
w=woord(c) // het woord met deze nummer in het woordenboek

output w

while (nog een code om te lezen)
{
    c= volgende code
    if (c in woordenboek)
        volgend_woord=woord(c)
    else
        volgend_woord= [w,w[0]]
        // de string w met het teken w[0] achteraan

    voeg [w,volgend_woord[0]] toe aan het woordenboek
        en geef het volgende nummer

    w=volgend_woord
    output w
}
```

Ook hier helpt een voorbeeld om te zien hoe het werkt:

<u>97 110 256 258 115</u>		
woord= a	tekst=a	nog voor de lus

97 110 256 258 115

woord= a    volgend woord= n  
woordenboek: 256=an    tekst=an

97 110 256 258 115

woord= n    volgend woord= an  
woordenboek: 257=na    tekst=anan

97 110 256 258 115

woord= an    volgend woord= ana (258 nog niet in woordenboek!)  
woordenboek: 258=ana    tekst=ananana

97 110 256 258 115

woord= ana    volgend woord= s  
woordenboek: 259=anas    tekst=anananas

Over het feit of LZW of LZ77 beter comprimeren vind je verschillende uitspraken. Daarbij zijn zeker vooral de precieze implementaties belangrijk en natuurlijk de toepassingen. Maar één voordeel heeft LZW zeker: het coderen kan op een efficiëntere manier gebeuren dan voor LZ77 omdat het dure opzoeken van een langst mogelijke deelstring niet moet gebeuren. Het zoeken of een deelstring al in het woordenboek zit, kan op een veel efficiëntere manier gebeuren (bv. met tries – bomen waar de bogen labels hebben zoals in de Huffman bomen).

### Oefening 105 *Codeer de teksten*

- `swiss_miss_missing`
- `een_beenen_been` (*met spelfout werkt het iets beter ...*)
- `ananas-kanaal`

*met LZW.*

**Oefening 106** *Decodeer de volgende met LZW gecodeerde tekst. De nieuwe woorden in het woordenboek beginnen met nummer 256. Geef ook de woorden met code ten minste 256 in het woordenboek. Wees niet verrast als de tekst niet echt zinnig is. . . Lettertakens worden hier als tekens gegeven en niet als hun code (om de ASCII code niet te hoeven opzoeken).*

a 256 b c a b 257 260

**Oefening 107** *Welke datastructuren zou je gebruiken om jouw woordenboek voor het LZW algoritme bij te houden? Bespreek het woordenboek dat je bij het coderen gebruikt en het woordenboek dat je bij het decoderen gebruikt apart.*

## 5.4 De Burrows-Wheeler transformatie.

De Burrows-Wheeler transformatie (Michael Burrows en David Wheeler) werd in een artikel van 1994 gepubliceerd, maar gaat terug op onderzoek dat Wheeler al in 1983 deed en toen (nog) niet publiceerde.

De Burrows-Wheeler transformatie is de basis van het programma *bzip2* dat in de meeste gevallen heel goede resultaten oplevert (beter dan bv. *gzip* of *zip*). Maar voor ons is het vooral interessant omdat het een fundamenteel nieuw idee bevat: waarom niet de tekst wijzigen (transformeren) zodat hij achteraf beter gecomprimeerd kan worden? Natuurlijk is het belangrijk dat de wijziging achteraf ongedaan gemaakt kan worden om de originele tekst terug te krijgen.

Wij zullen het algoritme uitleggen aan de voorbeeldtekst **boom\_gaat** (waarbij wij Lars Boom eens met een kleine b schrijven om geen verwarring te hebben). De Burrows-Wheeler transformatie gebruikt de tekst niet sequentieel, maar begint al met heel grote stukken van de tekst. Misschien gebruikt het de hele tekst, maar ten minste voldoende grote stukken (blokken). Hoe groter de blokken zijn hoe beter de compressie – maar vooral de compressie is ook duurder voor grotere blokken. In onze kleine voorbeelden zal de tekst natuurlijk helemaal in één blok getransformeerd kunnen worden.

Eerst zullen wij de tekst zo interpreteren alsof het geen string is, maar een cykel van karakters, dus alsof na het laatste letterteken terug het eerste letterteken komt. Stel dat de lengte van de tekst  $n$  is. Dan kan je van deze *cyclische tekst*  $n$  strings van lengte  $n$  krijgen door op de  $n$  verschillende posities te beginnen. In ons voorbeeld is  $n = 9$  en wij krijgen

b	o	o	m	_	g	a	a	t
o	o	m	_	g	a	a	t	b
o	m	_	g	a	a	t	b	o
m	_	g	a	a	t	b	o	o
_	g	a	a	t	b	o	o	m
g	a	a	t	b	o	o	m	_
a	a	t	b	o	o	m	_	g
a	t	b	o	o	m	_	g	a
t	b	o	o	m	_	g	a	a

Wij zien dat in elke rij en elke kolom de hele tekst staat. Als jullie het echt willen implementeren zouden jullie natuurlijk het best niet echt  $n$  kopieën van de tekst gebruiken, maar een enkele kopie en de verschillende strings gewoon voorstellen als pointers naar het begin karakter!

Nu sorteren wij de rijen. Daarbij gebruiken wij de waarden van de bytes, maar omdat onze voorbeelden voor compressie met gewone teksten met kleine lettertekens werken, komt dat hier overeen met de lexicografische volgorde. Dan krijgen wij

a <sub>1</sub>	a <sub>2</sub>	t	b	o <sub>1</sub>	o <sub>2</sub>	m	_	g
a <sub>2</sub>	t	b	o <sub>1</sub>	o <sub>2</sub>	m	_	g	a <sub>1</sub>
b	o <sub>1</sub>	o <sub>2</sub>	m	_	g	a <sub>1</sub>	a <sub>2</sub>	t
g	a <sub>1</sub>	a <sub>2</sub>	t	b	o <sub>1</sub>	o <sub>2</sub>	m	_
m	_	g	a <sub>1</sub>	a <sub>2</sub>	t	b	o <sub>1</sub>	o <sub>2</sub>
o <sub>2</sub>	m	_	g	a <sub>1</sub>	a <sub>2</sub>	t	b	o <sub>1</sub>
o <sub>1</sub>	o <sub>2</sub>	m	_	g	a <sub>1</sub>	a <sub>2</sub>	t	b
t	b	o <sub>1</sub>	o <sub>2</sub>	m	_	g	a <sub>1</sub>	a <sub>2</sub>
_	g	a <sub>1</sub>	a <sub>2</sub>	t	b	o <sub>1</sub>	o <sub>2</sub>	m

Daarbij hebben wij een index bij de o en bij de a geschreven om te kunnen zien *welk van de identieke lettertekens (bytes) het is*. Maar dat is alleen maar ter illustratie en niet deel van de methode.

Wat wij nu **eerst** stellen is dat alle rijen van deze matrix verschillend zijn. Als dat niet zo was, kan je dat gemakkelijk zien omdat dan in de gesorteerde matrix identieke rijen achter elkaar moeten staan. In dit geval is de hele tekst een herhaling van een deelstring en je kan het best de deelstring opslaan



(misschien nog gecomprimeerd) en hoe vaak je hem moet herhalen. Een andere oplossing zou zijn dat je op het einde van de originele string een teken toevoegt dat nog niet in de string zit en na de inverse transformatie opnieuw verwijderd moet worden. Dat kan bv. een teken zijn dat kleiner is dan alle tekens in de string (als dat kan). Op deze manier kan je op voorhand garanderen dat alle rijen verschillend zijn!

De getransformeerde tekst is nu een kolom van deze matrix. Voor compressie doeleinden zou het natuurlijk ideaal zijn de eerste kolom te gebruiken – in een lange tekst zouden daar in het begin alleen maar a's staan, dan alleen maar b's, etc – je zou dus alleen maar de hoeveelheden moeten bijhouden. Jammer genoeg is er geen manier om uit deze kolom de originele (cyclische) tekst te reconstrueren. Maar als het om een *echte* en niet om een toevallige tekst gaat, heeft ook de laatste kolom nog voordelen van het sorteren: ook hier is de verdeling van de lettertekens niet toevallig omdat de lettertekens de prefixen van de lettertekens in de eerste kolom zijn. En wat verrassend is: terwijl je uit de eerste kolom de cyclische tekst niet terug kan krijgen (zonder dure extra informatie op te slaan) kan je dat uit de laatste kolom wel! Je moet alleen nog bijhouden waar het eerste letterteken  $s[0]$  van de string  $s[]$  staat om ook de niet cyclische string te reconstrueren.

Een voordeel is dat als wij de laatste kolom hebben, wij de eerste kolom gemakkelijk kunnen verkrijgen: het zijn dezelfde lettertekens als in de laatste kolom en wij moeten die gewoon sorteren! Daarbij houden wij natuurlijk geen rekening met de indices (tenslotte zijn die er niet echt), maar stel voor het moment dat je ook de indices kent. Wij hebben dus

$a_1$	.....	$g$
$a_2$	.....	$a_1$
$b$	.....	$t$
$g$	.....	—
$m$	.....	$o_2$
$o_2$	.....	$o_1$
$o_1$	.....	$b$
$t$	.....	$a_2$
—	.....	$m$

Het transformatiealgoritme kan dus als volgt geschetst worden:

- interpreteer de tekst  $s[0], s[1], \dots, s[n-1]$  als cyclisch en vorm alle  $n$  strings  $s_i[]$  met  $0 \leq i < n$  die je krijgt door op positie  $i$  te beginnen.

- sorteer deze  $n$  strings in stijgende lexicografische volgorde. Wij schrijven  $p_s(i)$  voor de string op positie  $i$  na het sorteren –  $p_s(0)$  is dus de kleinste string.
- output het nummer  $i$  van de rij die met  $s[0]$  eindigt – dus de rij  $i$  met  $p_s(i)[n-1] = s[0]$ . Dat is de rij waar het letterteken waarmee de originele string begint op de laatste positie staat. Output dan een blank.
- output de lettertekens  $p_s(0)[n-1], \dots, p_s(n-1)[n-1]$  op de laatste posities in deze volgorde.

Het eerste teken van de originele tekst staat in ons geval in rij 6 als wij met 0 beginnen te tellen. De Burrows-Wheeler getransformeerde van **boom\_gaat** zou dus bv. **(6\_gat\_oobam)** zijn als wij vastleggen dat het getal eerst komt en de tekst na de eerste blank staat.

Maar nu is het gemakkelijk de cyclische tekst opnieuw op te bouwen. Wij moeten voor elk letterteken alleen het volgende letterteken kennen – maar na het letterteken dat in de laatste kolom in rij  $i$  staat komt het letterteken dat in de eerste kolom in rij  $i$  staat – en dat kennen wij!

Als wij dus met het eerste teken van de tekst beginnen – dus het teken nummer 6 van **gat\_oobam** dan is dat **b** en in de eerste kolom staat op positie 6 **o**<sub>1</sub> – het volgende teken is dus **o**<sub>1</sub>. In de laatste kolom staat **o**<sub>1</sub> in rij 5 en in rij 5 in de eerste kolom staat **o**<sub>2</sub>. Dat vinden wij terug in rij 4 in de laatste kolom, etc. . . Als wij zo doorgaan krijgen wij de oorspronkelijke tekst **boom\_gaat**.

Er blijft maar één probleem: toen wij opgezocht hebben welk van de o's aan de rechterkant de **o** is die naar de **b** komt, hebben wij de indices gebruikt – maar die zijn er in het echt niet – daar zijn het alleen maar a's en o's!

Als jullie naar de volgorde kijken waarop de indices opduiken dan is het in de eerste kolom eerst **a**<sub>1</sub> en dan **a**<sub>2</sub> – en voor de tweede kolom ook. Voor de o's is het eerst **o**<sub>2</sub> en dan **o**<sub>1</sub> in de eerste kolom – en in de tweede kolom ook. En dat is inderdaad geen toeval: wij kunnen identieke lettertekens in de eerste en laatste kolom aan elkaar toekennen in de volgorde waarin ze opduiken!

Om dat te bewijzen stel dat de string  $s[] = s[0] \dots s[n-1]$  is en dat  $s[i] = s[j]$  waarbij  $s[i]$  in de gesorteerde matrix in de eerste kolom voor  $s[j]$  komt. Wij zouden nu graag bewijzen dat dan ook in de laatste kolom  $s[i]$  voor  $s[j]$  komt, maar jammer genoeg is dat niet altijd zo:

$$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \begin{pmatrix} a_2 \\ a_1 \end{pmatrix} \longrightarrow \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$$

$$\begin{pmatrix} a_1 \\ b_1 \\ a_2 \\ b_2 \end{pmatrix} \begin{pmatrix} b_1 \\ a_2 \\ b_2 \\ a_1 \end{pmatrix} \begin{pmatrix} a_2 \\ b_2 \\ a_1 \\ b_1 \end{pmatrix} \begin{pmatrix} b_2 \\ a_1 \\ b_1 \\ a_2 \end{pmatrix} \longrightarrow \begin{pmatrix} a_1 \\ a_2 \\ b_1 \\ b_2 \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \\ a_2 \\ a_1 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ b_1 \\ b_2 \end{pmatrix} \begin{pmatrix} b_2 \\ b_1 \\ a_2 \\ a_1 \end{pmatrix}$$

Maar als wij stellen dat de rijen allemaal verschillend zijn, kunnen wij gemakkelijk bewijzen dat de volgorden in de eerste en laatste kolom dezelfde zijn:

Stel dat  $s[i] = s[j]$  en dat  $s[i]$  in de gesorteerde matrix in de eerste kolom voor  $s[j]$  komt – dus

$$(s[i], s[i+1] \dots s[n-1]s[0] \dots s[i-1]) < (s[j], s[j+1] \dots s[n-1]s[0] \dots s[j-1]).$$

Omdat  $s[i] = s[j]$  geldt dus

$$(s[i+1] \dots s[n-1]s[0] \dots s[i-1]) < (s[j+1] \dots s[n-1]s[0] \dots s[j-1]) \text{ en dus ook}$$

$$(s[i+1] \dots s[n-1]s[0] \dots s[i-1]s[i]) < (s[j+1] \dots s[n-1]s[0] \dots s[j-1]s[j]).$$

Maar dat betekent precies dat in de laatste kolom de rij waarin  $s[i]$  staat voor de rij komt waarin  $s[j]$  staat.

Maar precies dat wouden wij bewijzen en dus is onze methode om lettertekens met elkaar te identificeren (indices toe te kennen) onder deze omstandigheden juist.

Het terugtransformatiealgoritme kan dus als volgt geschetst worden:

Stel dat de code  $(p, c[0], \dots, c[n-1])$  is.

- kopieer  $c[0], \dots, c[n-1]$  en sorteer de tekens. Noem de nieuwe string  $e[]$ . In  $e[i]$  staat dus het letterteken dat in de originele tekst na het teken  $c[i]$  komt.
- voor elk letterteken  $x$  dat  $k$  keer in  $c[]$  en dus ook  $k$  keer in  $e[]$  zit, ken indices  $1, \dots, k$  toe aan de posities van  $x$  in  $c[]$  in de volgorde van de string. Doe hetzelfde voor  $e[]$ . Wij schrijven  $x_k$  voor het letterteken  $x$  met index  $k$ . Elk geïndexeerd letterteken zit dus één keer in  $c[]$  en één keer in  $e[]$ .
- begin met  $x_j = c[p]$  dan herhaal  $n$  keer:

- output het letterteken  $x$  (zonder index).
- bepaal het nummer  $i$  zodat  $c[i] = x_j$ .
- kies  $x_j = e[i]$  als volgend geïndexeerd letterteken.

Jullie zullen op het internet ook beschrijvingen vinden waar jullie de lus niet  $n$  keer moeten herhalen, maar *totdat je terug bent bij het eerste letterteken*. Maar wij zullen zien dat het voordelen heeft het zo te doen...

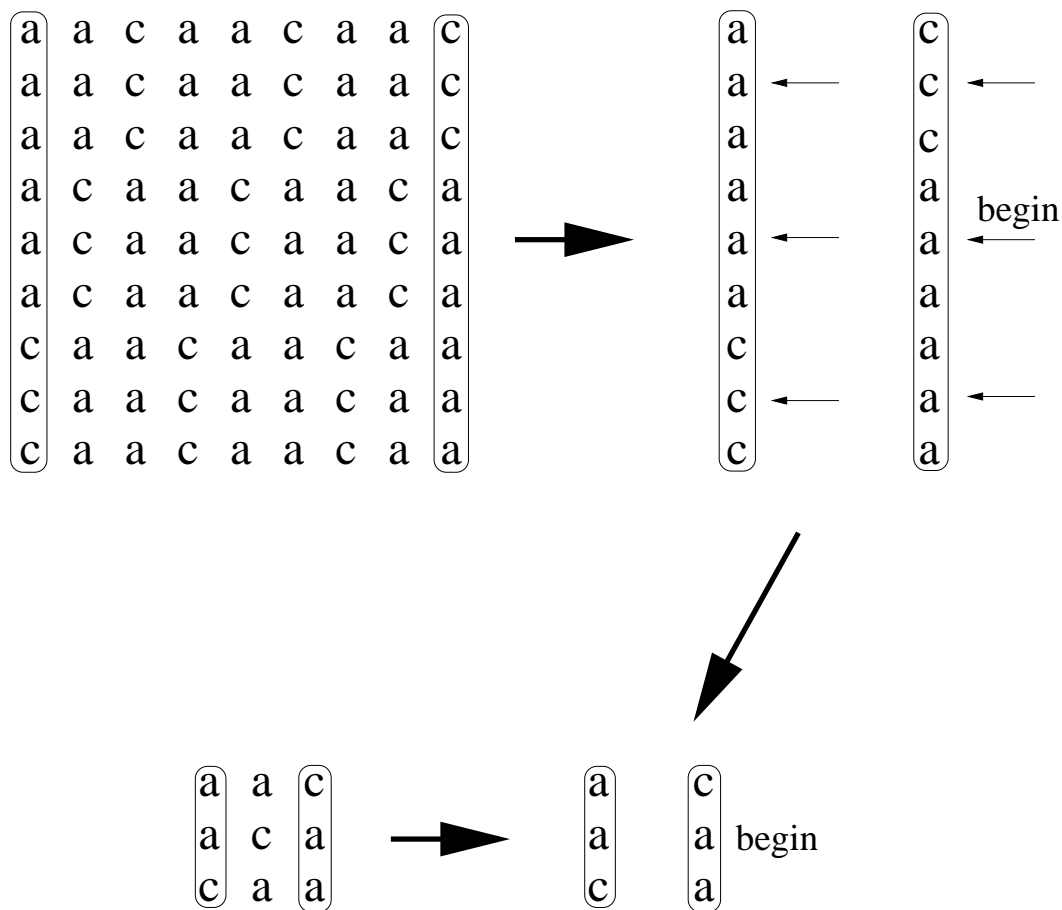
Wat zou er gebeuren als wij het algoritme zouden toepassen op een string die een herhaling van deelstrings is – dus waar twee of meerdere rijen identiek zijn?

Stel dat de string  $t[0], \dots, t[k-1]$  niet tot identieke rijen leidt maar de hele string  $s[]$  wel, omdat hij uit  $l$  herhalingen van  $t[]$  bestaat. Dus  $s[] = t[0], \dots, t[k-1], t[0], \dots, t[k-1], \dots, t[0], \dots, t[k-1]$ . Als wij op deze string het algoritme toepassen, hebben wij elke rij precies  $l$  keer achter elkaar. De hele matrix is dus een opeenvolging van  $k$  blokken waarvan elke blok  $l$  identieke rijen bevat. Rij in verschillende blokken zijn verschillend. Maar voor de transformatie doet dat er natuurlijk niet toe – wij kunnen gewoon de getransformeerde tekst construeren.

Nu kijken wij naar het toekennen van indices in  $c[]$  (en analoog in  $e[]$ ). Tijdens de terugtransformatie zal een teken in  $e[]$  of  $c[]$  dat op positie  $i < l$  in zijn eigen blok zit zeker een index  $j$  krijgen zodat  $i = j(\bmod l)$ : in zijn blok is het zeker het  $i$ -de en in de blokken daarvoor was het identieke letterteken ofwel niet op de eerste (of laatste) positie aanwezig ofwel  $l$  keer – dus zal de index  $j$  in elk geval aan  $i = j(\bmod l)$  voldoen.

Maar als ons beginteken een index  $b$  heeft dan zien wij dat in het hele algoritme **nooit** een letterteken met een index  $i$  gebruikt wordt die niet aan  $i(\bmod l) = b(\bmod l)$  voldoet. Het algoritme gebruikt dus alleen maar één rij van elk blok. Maar als wij naar de eerste en laatste lettertekens van alleen deze rijen kijken dan zien wij dat dit dezelfde zijn als bij het toepassen van het algoritme alleen op  $t[0], \dots, t[k-1]$ . Een voorbeeld ervan zien wij in Figuur 5.4. Inderdaad zullen wij dus na  $k$  stappen terug zijn met ons beginteken – maar omdat wij even veel stappen doen als er lettertekens in de getransformeerde tekst zitten – dus  $l * k$ , zullen wij doorgaan en de tekst  $t[0], \dots, t[k-1]$   $l$  keer herhalen. De originele tekst wordt dus inderdaad ook in dit geval juist teruggetransformeerd.

Het lijkt natuurlijk niet echt efficiënt als wij het grootste deel van de getransformeerde tekst helemaal niet gebruiken – en ook niet nodig hebben. Maar wij zullen zien dat in dit geval de volgende stappen ook bijzonder goed comprimeren.



Figuur 18:

Wat hebben wij tot nu toe gewonnen? De getransformeerde tekst heeft even veel lettertekens. En als wij nu het Huffman algoritme zouden toepassen, hebben wij nog altijd niets gewonnen: de relatieve aantallen lettertekens zijn dezelfde als in de originele tekst – de code zou dus even lang zijn (in feite zelfs iets langer omdat het begin van de tekst ook opgeslaan moet worden).

### Wat is dan het voordeel?

De volgorde van de strings is lexicografisch. De meeste invloed daarbij heeft natuurlijk het eerste letterteken in elke rij – daar hebben wij een mooie ordening. Maar het laatste letterteken is een prefix van het eerste. In een toevallige tekst betekent dat niet veel omdat elk teken dezelfde kans heeft een prefix te zijn. Maar wij weten al dat je toevallige data niet kan comprimeren. In *echte* teksten (en ook andere bestanden) zijn dergelijke prefixen niet toevallig. De strings die bv. met *aa* beginnen, zullen allemaal achter

elkaar staan. Maar het is zeker niet zo dat elk letterteken dezelfde kans heeft een prefix van *aa* te zijn. Als jullie bv. in het groene boekje kijken dan zijn er veel woorden die met **baa** of **laa** beginnen, duidelijk minder met **faa** maar geen enkel woord begint met **caa**. Dat zegt natuurlijk niet alles, omdat natuurlijk de prefixen niet altijd op het begin van een woord staan, maar het geeft al een idee. En als jullie naar langere deelstrings kijken, wordt dat nog duidelijker (bv. de prefixen van *aan* of nog langere strings...). Wij kunnen dus hopen dat er in de rijen waar *aa* de volgorde bepaalt op het einde veel b's of l's bij elkaar staan. Nog duidelijker wordt het als jullie naar langere strings kijken. Alle strings die met *eze* beginnen, zullen achter elkaar staan – en de meeste ervan zullen zeker op de laatste positie een **d** hebben. In onze code zullen dus veel d's op elkaar volgen of ten minste heel dicht bij elkaar staan. De letters zijn dus beter gegroepeerd. Wij hebben dus een compressiealgoritme nodig waarvoor dat een voordeel is – adaptive Huffman coding zou bv. al beter presteren!

Maar Burrows en Wheeler stellen voor deze eigenschap van de getransformeerde tekst als volgt te gebruiken:

### ***De move to front methode***

Je legt vast hoe je jouw alfabet met getallen codeert – je houdt de tekens bij in een lijst en elk letterteken  $x$  krijgt als code  $c(x)$  zijn positie in de lijst. In de realiteit zouden dat bv. alle ASCII codes zijn – omdat je niet op voorhand weet welke lettertekens opduiken en welke niet moet je een **vaste** lijst hebben met **alle** mogelijke lettertekens! In onze voorbeelden willen wij natuurlijk niet met zo'n lange lijst werken, wij kiezen dus een kortere, maar het is belangrijk om te onthouden dat dat alleen voor de voorbeelden is en in het algemeen niet kan!

**Voorbeeld:** De lijst: (6,a,b,c,g,m,o,t,z,-). Dus is  $c(6) = 0$ ,  $c(a) = 1, \dots, c(-) = 9$  etc.

Als jullie nu een string coderen, dan schrijven jullie voor elk letterteken  $x$  zijn code  $c(x)$  **en** plaatsen dit letterteken dan in het begin van de lijst – het heeft **achteraf** dus code 0.

Als jullie met de lijst hierboven dus **6\_gat\_oobam** willen coderen dan werkt dat als volgt:

de lijst is: (6,a,b,c,g,m,o,t,z,-) en **6\_gat\_oobam** moet gecodeerd worden.

lees:**6**: 0 schrijven, lijst nu (6,a,b,c,g,m,o,t,z,-)

lees:**-**: 9 schrijven, lijst nu (-,6,a,b,c,g,m,o,t,z)

lees:**g**: 5 schrijven, lijst nu (g,-,6,a,b,c,m,o,t,z,-)

lees:**a**: 3 schrijven, lijst nu (a,g,\_,6,b,c,m,o,t,z,)  
 lees:**t**: 8 schrijven, lijst nu (t,a,g,\_,6,b,c,m,o,z,)  
 lees:**\_**: 3 schrijven, lijst nu (\_,t,a,g,6,b,c,m,o,z)  
 lees:**o**: 8 schrijven, lijst nu (o,\_,t,a,g,6,b,c,m,z)  
 lees:**o**: 0 schrijven, lijst nu (o,\_,t,a,g,6,b,c,m,z)  
 lees:**b**: 6 schrijven, lijst nu (b,o,\_,t,a,g,6,c,m,z)  
 lees:**a**: 4 schrijven, lijst nu (a,b,o,\_,t,g,6,c,m,z)  
 lees:**m**: 8 schrijven, lijst nu (m,a,b,o,\_,t,g,6,c,z)

Dit toont wel het principe, maar het kleine voorbeeld is niet voldoende om aan te tonen dat het nuttig is. Maar omdat in de Burrows-Wheeler getransformeerde tekst heel vaak herhalingen van letters voorkomen – of ten hoogste met kleine onderbrekingen, zullen in de code zeer veel nullen en kleine getallen opduiken.

Als je het decodeert moet je natuurlijk met dezelfde lijst beginnen en die dan ook op dezelfde manier veranderen:

de lijst is: (6,a,b,c,g,m,o,t,z,\_) en **0,9,5,3,8,3,8,0,6,4,8** moet gedecodeerd worden.

lees:**0**: 6 schrijven, lijst nu (6,a,b,c,g,m,o,t,z,\_)  
 lees:**9**: \_ schrijven, lijst nu (\_,6,a,b,c,g,m,o,t,z)  
 lees:**5**: g schrijven, lijst nu (g,\_,6,a,b,c,m,o,t,z,)  
 lees:**3**: a schrijven, lijst nu (a,g,\_,6,b,c,m,o,t,z,)  
 lees:**8**: t schrijven, lijst nu (t,a,g,\_,6,b,c,m,o,z,)  
 lees:**3**: \_ schrijven, lijst nu (\_,t,a,g,6,b,c,m,o,z)  
 lees:**8**: o schrijven, lijst nu (o,\_,t,a,g,6,b,c,m,z)  
 lees:**0**: o schrijven, lijst nu (o,\_,t,a,g,6,b,c,m,z)  
 lees:**6**: b schrijven, lijst nu (b,o,\_,t,a,g,6,c,m,z)  
 lees:**4**: a schrijven, lijst nu (a,b,o,\_,t,g,6,c,m,z)  
 lees:**8**: m schrijven, lijst nu (m,a,b,o,\_,t,g,6,c,z)

Omdat er heel veel nullen en kleine getallen inzitten is de met de move to front methode gecodeerde Burrows-Wheeler getransformeerde tekst nu heel geschikt om (statistisch) Huffman gecodeerd te worden – en dat stellen Burrows en Wheeler inderdaad ook als één van de mogelijkheden voor.

**Oefening 108** *Codeer de tekst jan kan van alles door eerst de Burrows-Wheeler transformatie toe te passen en dan de move to front codering. De lijst voor de codering is (a,b,e,j,k,l,m,n,s,v,\_,0,1,2,3,4,5,6,7,8,9).*

**Oefening 109** *Decodeer de volgende code: (0,11,10,10,7,0,6,1,0,9,11,4,6,1)*

*De code is de move to front code van een Burrows wheeler getransformeerde tekst waarbij het move to front algoritme met de lijst (7,a,e,i,o,u,b,c,n,t,w,-) begonnen is. De 7 is hier natuurlijk een gewoon letterteken en geen getal!*

**Oefening 110** *Als wij de eerste kolom als getransformeerde tekst nemen, kunnen wij de oorspronkelijke tekst niet reconstrueren. Maar hoe zit het met de tweede kolom. Werk een algoritme voor de tweede kolom uit of geef een redenering waarom de tweede kolom als getransformeerde tekst niet geschikt is.*

**Oefening 111** *Zou het algoritme werken als je in plaats van de strings in stijgende volgorde te sorteren de strings in dalende volgorde zou sorteren? Zouden er wijzigingen nodig zijn? Kijk bv. naar de redenering dat de volgorde van lettertekens in de eerste en laatste kolom dezelfde zijn. Bewijs dat dat ook met dalende volgorde het geval is of geef een tegenvoorbeeld.*

**Oefening 112** *Geef een string  $x t_0 \dots t_{n-1}$  waarbij  $x$  de voorstelling van een getal  $x'$  is met  $0 \leq x' < n$  en  $t_0, \dots, t_{n-1}$  lettertekens. Maar deze string mag niet het resultaat van een Burrows-Wheeler transformatie kunnen zijn. Geef uitleg.*

**Oefening 113** *Werkt het volgende algoritme dat op het idee van de Burrows-Wheeler transformatie gebaseerd is juist? Bewijs dat of geef een tegenvoorbeeld.*

*Neem de tekst  $t = t_0, \dots, t_{n-1}$  en voeg het teken  $t_n = \infty$  toe. Daarbij betekent  $\infty$  dat het een teken is dat groter is dan welk teken ook in het alfabet. Schrijf dan  $s_i$  voor de string  $t_i, \dots, t_n$  en sorteer de strings  $s_0, \dots, s_n$  in lexicografische volgorde. Als nu de string  $s_k$  op positie  $j$  in deze volgorde staat dan is  $t_{k-1}$  het  $j$ -de teken van de getransformeerde tekst (met  $t_{-1} = t_n$ ).*

*De terugtransformatie begint met het teken  $\infty$  – maar zonder het te schrijven – en gaat dan door zoals in de les gezien totdat  $n$  tekens geschreven zijn.*

## **6 Parallele algoritmen**

*Parallele algoritmen* zijn algoritmen die meerdere processoren gebruiken die samenwerken om één probleem op te lossen. Het voordeel is duidelijk en het principe van parallelisme wordt al lang gebruikt: als bv. een huis wordt gebouwd dan zou (bijna) niemand proberen dat helemaal alleen te doen – het aantal werkuren is zo groot dat de mens die het huis wil bouwen misschien al



overleden is als het huis eindelijk klaar is. . . Daarom werken meerdere mensen tegelijk. Daardoor wordt het aantal werkuren niet kleiner (misschien zelfs groter) maar de tijd tussen begin van de bouw en het einde wordt veel kleiner. Dit is ook het doel van parallelle algoritmen. Het is duidelijk dat de totale verbruikte rekentijd niet kleiner kan worden (je kan de taken van de verschillende processoren ook na elkaar op één processor laten draaien) maar de *wachttijd* tussen het begin van de berekening en het einde kan wel duidelijk korter worden.

Maar het voorbeeld met ploegen die samen werken toont ook twee problemen: soms moet je goed op de volgorde letten (bv. de muren niet schilderen als er nog kabels in de muren gelegd moeten worden) en je kan ook niet arbitrair veel werkers gebruiken. Het helpt bv. niet als 10 personen proberen een enkele schakelaar voor het licht vast te maken. Het gemakkelijkst is het taken aan verschillende werknemers toe te kennen als de taken helemaal niets met elkaar te maken hebben (bv. de kelder schilderen en tegelijk pannen op het dak plaatsen).

Hetzelfde geldt voor parallelle algoritmen: als deeltaken niets met elkaar te maken hebben, is het veel gemakkelijker ze op verschillende processoren te laten draaien dan in gevallen waar de taken sterk van elkaar afhankelijk zijn.

Een probleem bij de beslissing welke onderwerpen voor dit deel van de les gekozen worden, was dat het heel moeilijk is om te voorspellen wat in de toekomst **echt** relevant zal zijn. En dat geldt vooral voor een les waar de klemtoon vooral op een realistische kijk op de zaak ligt. Rond 1980 leek het (voor een tijdje) dat de kloksnelheid van de computers niet meer zo snel groeide en in de tijd na 1980 werd verwacht dat *parallel computing* de toekomst zou zijn. Er werden veel verschillende modellen voor parallelle computers ontworpen. Verschillende architecturen bestonden zelfs als hardware – bv. de bekende Cray supercomputers (vanaf 1976). Maar de transputers waar interprocessorcommunicatie een belangrijk deel van het ontwerp was waren bijzonder interessant. Deze ontwerpen besteedden veel tijd en geld aan het ontwerp van een efficiënt netwerk voor de interprocessorcommunicatie.

Maar het “probleem” was vermoedelijk dat door de extreem grote aantallen van kleine computers er veel meer geld aan de ontwikkeling van CPU’s besteed kon worden die geoptimaliseerd zijn om zelfstandig te werken. Daardoor verdwenen de *echte* parallelle computers langzaam en in plaats daarvan ontstonden clusters van stand-alone computers die door netwerken met elkaar verbonden waren die – in vergelijking met de kloksnelheid – zeer traag zijn. Deze ontwikkeling kan je – natuurlijk – ook in de opleiding terugvinden, waarin vroeger een hele les *parallelle algoritmen* gegeven werd die later gereduceerd werd op een relatief klein deel van DA 3.

Intussen is het opnieuw zo dat de kloksnelheden niet meer zo snel groeien – maar de processoren worden door het feit dat per clockcycle meer en ingewikkeldere bewerkingen plaats kunnen vinden wel nog altijd sneller. Toch lijkt parallelisme opnieuw belangrijker te worden (wat in toekomst misschien ook tot een groter aandeel van parallelle algoritmen in de les of zelfs een eigen les kan leiden). Bijzonder interessant zijn daarbij de multicore computers die het idee van een shared memory verwezenlijken en ook in grote hoeveelheden verkocht worden. Er gebeurt ook veel onderzoek over de mogelijkheid de grafische kaarten – die massief parallel zijn – voor berekeningen te gebruiken. De vraag is dus in welke richting de ontwikkeling deze keer gaat. Parallelle algoritmen hebben het normaal niet over 2, 4 of 8 processoren, maar over zoiets als bv.  $O(n)$  als  $n$  de grootte van het probleem is. De vraag is dus of er een ontwikkeling in de richting van multicore processoren met duidelijk meer cores zal zijn of of de ontwikkeling zal stoppen en opnieuw in een andere richting gaan. . .

Misschien zeggen jullie nu dat er toch heel veel clusters zijn waarop *parallel computing* wordt toegepast – bv. voor voorspellingen van het weer en ook *scientific computing*. Maar in die gevallen wordt eerder iets toegepast dat *distributed computing* genoemd wordt. Het verschil is dat in *distributed computing* de communicatie tussen de verschillende delen *relatief* klein is. In een cluster gebeurt de communicatie via een intern netwerk, maar in sommige gevallen van distributed computing gebeurt de communicatie van de computers die de deeloplossingen berekenen zelfs over het internet. De communicatie tussen verschillende rekennodes is dus **duidelijk** trager dan de communicatie van een CPU met zijn lokaal geheugen. Voor dergelijke algoritmen is de precieze vorm van het netwerk dan ook minder belangrijk omdat de communicatie tussen de delen klein is.

De klemtoon van deze les ligt wel op realistische toepassingen eerder dan theoretische resultaten, maar toch zullen jullie in dit deel van allebei iets zien. . .

**Notatie:** Wij schrijven  $t_p(n)$  voor de tijd die een gegeven parallel algoritme dat  $p$  processoren voor een taak met invoerlengte  $n$  gebruikt ten hoogste nodig heeft. Dat is dus niet de som van de tijden van de enkele processoren, maar gewoon de tijd die verstreken is tussen het begin van de berekeningen en het einde. Wij schrijven  $t^s(n)$  voor de tijd die het snelste (gekende) sequentiële algoritme voor deze input nodig heeft. Het is onmiddellijk duidelijk dat als een parallel algoritme in tijd  $O(t_p(n))$  draait dat dan  $O(p * t_p(n))$  een bovengrens is voor  $t^s(n)$  omdat door achter elkaar uitvoeren van de stappen van de  $p$  processoren het parallelle algoritme ook serieel gesimuleerd kan worden, maar de simulatie kan natuurlijk wel een zekere (constante) overhead

vragen.

Het product van de tijd en het aantal processoren wordt soms ook de kost  $K(n)$  van het algoritme genoemd. Wij weten door de mogelijkheid van simulatie dat dit product (op een constante na) optimaal voor seriële algoritmen is en wij noemen een parallel algoritme kost optimaal als  $K(n) = O(t(n))$  waarbij  $t(n)$  de tijd voor het beste seriële algoritme is.

Wij zullen nog een ander concept zien: de werk complexiteit van een algoritme. Als de scheduler op een parallelle computer slim genoeg is dan kan hij in sommige gevallen herkennen dat sommige processoren na de eerste stappen niet meer gebruikt worden. Dat is niet in elk geval gemakkelijk, maar in principe kan de programmeur zelf in sommige gevallen misschien een `exit()` in de code voor een processor  $i$  schrijven. Dan kunnen deze processoren voor andere doeleinden gebruikt worden. Het is dus niet in elk geval zinvol zo te rekenen alsof alle ooit gebruikte processoren de hele tijd bezig waren als het om de kost van een algoritme gaat. De *kost* die de tijd van de processoren alleen telt als deze processoren ook stappen voor het algoritme uitvoeren wordt de *werk complexiteit* genoemd. Of precies: het werk  $W(n)$  dat een parallel algoritme gebruikt voor een invoer met lengte  $n$  is gedefinieerd als

$$W(n) = \sum_{i=1}^p t(i, n)$$

waarbij processoren  $1, \dots, p$  aan de taak werken en  $t(i, n)$  het aantal stappen (de tijd) is dat processor  $i$  maximaal doet voor een invoer met lengte  $n$ . Zoals het hier gedefinieerd is, is het niet noodzakelijk zo dat de tijd die voor verschillende processoren gerekend wordt altijd voor dezelfde invoer maximaal is. Maar in de gevallen die wij hier zullen zien is dat wel zo...

Het is duidelijk dat een algoritme dat werk  $W(n)$  vraagt op een seriële machine in tijd  $O(W(n))$  gesimuleerd kan worden en analoog met *kost optimaal* definiëren wij dat een parallel algoritme werk optimaal is als  $W(n) = O(t(n))$  waarbij  $t(n)$  de tijd voor het beste seriële algoritme is.

## 6.1 Branch and bound met verdeelde processoren

Als het om problemen gaat met een heel goedkope asymptotische complexiteit (bv.  $O(n)$  of  $O(n \log n)$  met een niet te grote constante voor de  $O()$ ) dan moet je al extreem grote hoeveelheden van data hebben om ervoor te zorgen dat meer dan één processor nodig is om het probleem in aanvaardbare tijd af te werken. Bovendien is de bottleneck dan misschien zelfs eerder het lezen en schrijven van de data. Dat gebeurt dan het best ook verdeeld.

Maar de duurste problemen die jullie in Datastructuren en Algoritmen 2 hebben gezien, zijn NP-complete problemen waarvoor de voor de hand liggende oplossing vaak branch and bound algoritmen zijn (hoewel andere aanzetten

– zoals bv. met integer programming – soms heel succesvol zijn). Hier zullen wij het erover hebben hoe je een branch and bound algoritme met een eenvoudig trucje efficiënt op verschillende machines kan verdelen. In een abstracte vorm kan je zo'n branch and bound algoritme ongeveer als volgt schrijven (waarbij in dit voorbeeld een minimum wordt gezocht):

**Algoritme 17** (*Branch and bound*)

```

beste_waarde= +oneindig; // of een gemakkelijk berekenbare
                        // bovengrens

b_and_b(configuratie)
// test de volgende configuratie in de recursieboom
{

if (is_endconfiguratie(configuratie)) // alle keuzes gemaakt
{ if (waarde(configuratie)<beste_waarde)
    beste_waarde=waarde(configuratie);
  return;
}

else
{
  benedengrens=bepaal_benedengrens(configuratie);
  // dat is de look-ahead van branch and bound die vaak moeilijk
  // is om te vinden maar belangrijk voor de performantie. Hier
  // wordt een benedengrens voor de waarde van een mogelijke
  // eindconfiguratie bepaald die uit deze configuratie voorkomt.

  if (benedengrens<beste_waarde)
  // verbetering nog mogelijk
  {
    for (elke mogelijke directe opvolgconfiguratie)
      b_and_b(opvolgconfiguratie);
  }
}
} // end else
return;
} // end b_and_b

```

Algoritme 17 wordt dan opgestart met de beginconfiguratie (de wortel van

de recursieboom) en op het einde staat de beste waarde in **beste\_waarde**. Om deze abstracte beschrijving iets beter te verstaan, kunnen wij bv. naar het inpakprobleem kijken, waar gewichten  $g_1 \leq 1, \dots, g_n \leq 1$  op een manier op vrachtwagens met capaciteit 1 geplaatst moeten worden dat het totale aantal vrachtwagens minimaal is. De *beginconfiguratie* is dan bv. de situatie waar alleen gewicht  $g_1$  op vrachtwagen 1 geplaatst is (dat mogen wij zeker stellen), een opvolgconfiguratie van een configuratie waar  $g_1, \dots, g_k$  geplaatst zijn is een configuratie waar  $g_1, \dots, g_{k+1}$  geplaatst zijn en een eindconfiguratie is een configuratie waar  $g_1, \dots, g_n$  geplaatst zijn. Een gemakkelijk berekenbare bovengrens is bv.  $n$ . De waarde van een configuratie is het aantal gebruikte vrachtwagens en **bepaal\_benedengrens(configuratie)** is bv. de waarde van **configuratie** plus de som van de nog niet geplaatste gewichten die op geen vrachtwagen meer passen die al bestaat (een slechte benedengrens, maar het gaat hier alleen om het principe).

Als jullie naar de verschillende takken van de recursieboom kijken dan valt op dat de enige samenhang de waarde **beste\_waarde** is.

Eén aanzet zou dus bv. zijn alle configuraties  $c_1, \dots, c_k$  op een zeker niveau – dat is de afstand  $x$  van de wortel – van de recursieboom op te slaan en die in parallel als beginconfiguraties te gebruiken. Elke keer dat een betere waarde voor **beste\_waarde** wordt gevonden, wordt die aan alle computers doorgestuurd. In de meeste gevallen zal het zo zijn als bij het inpakprobleem – dus dat **beste\_waarde** in vergelijking met het aantal toppen in de recursieboom niet vaak verbeterd kan worden. Bij het inpakprobleem kan **beste\_waarde** ten hoogste  $n - 1$  keer verbeterd worden terwijl de recursieboom exponentieel groot kan zijn! De communicatie tussen computers die verschillende takken afwerken is dus **miniem** en het is zelfs geen probleem de betere waarden via het internet te sturen!

Dat ziet er al niet slecht uit – maar er zijn twee problemen:

- a.) In de meeste gevallen zullen de takken allesbehalve even groot zijn. Sommige takken zullen zeer snel als slecht herkenbaar zijn (de bounding criteria zijn in deze delen heel efficiënt) en dus nauwelijks werk vragen en sommige takken zullen bijna de hele tijd vragen. Als wij bv. 100 computers ter beschikking hebben en de taak op deze manier in honderd delen splitsen dan is het mogelijk dat 99 delen na weinige seconden gedaan zijn en dat één van de taken weken duurt...
- b.) De verschillende startconfiguraties moeten berekend en op een manier opgeslagen worden dat de programma's vanuit dit punt kunnen beginnen. Dat is niet echt moeilijk, maar wij zullen zien dat het vaak overbodig werk is.

Probleem a.) kan gedeeltelijk opgelost worden door de afstand van de wortel van de recursieboom duidelijk groter te kiezen. Daardoor krijgen wij veel meer delen dan wij computers hebben en hoewel ook deze delen zeker verschillend veel tijd nodig hebben, kunnen wij hopen dat door het feit dat elke computer meerdere delen moet afwerken het verschil tussen de tijd die de computer nodig heeft die het langst bezig is (dat is onze wachttijd) en die het minst bezig is niet te groot wordt. Ten slotte kunnen computers die heel goedkope delen opstarten ook vroeger met het volgende deel beginnen en werken daardoor misschien meer delen af. Een probleem is natuurlijk dat wij misschien extreem veel delen hebben...

Branch and bound algoritmen zijn natuurlijk alleen duur als de bomen heel groot zijn. Dat betekent meestal niet dat de bomen heel diep zijn, maar dat er een grote vertakking is. Of met andere woorden: dat er op de lage niveaus (dus dicht bij de wortel) in vergelijking met de hele boom verwaarloosbaar weinig toppen zijn.

Om probleem b.) op te lossen zullen wij het gedeelte tot het gekozen niveau gewoon elke keer opnieuw doorlopen. Omdat het relatief klein is, verliezen wij op deze manier niet veel tijd. Het aantal gevallen mag natuurlijk niet te groot zijn! Om het aantal gevallen klein te houden en aan de andere kant een goede verdeling te krijgen, kiezen wij het level waar wij splitsen relatief diep in de boom (dus veel verschillende toppen om te splitsen), maar vatten wij toppen uit verschillende delen van de zoekboom samen. Of precies:

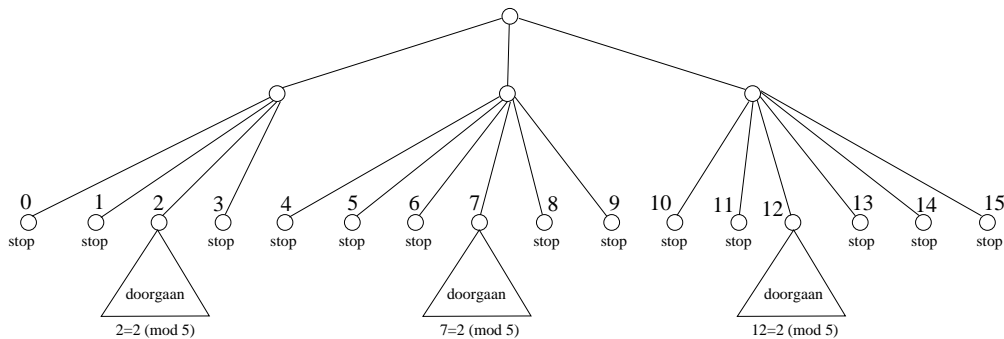
- Kies een diepte  $d$  om te splitsen en een aantal  $k$  van delen waarin het probleem gesplitst moet worden. Het aantal toppen van de recursieboom op diepte  $d$  moet zeer veel groter zijn dan  $k$ . De verschillende delen waarin de recursieboom opgesplitst wordt, kunnen nu beschreven worden door de indices  $0, 1, \dots, k - 1$ .
- Houd een teller bij die – beginnend met 0 – nummers toekent aan de toppen op diepte  $d$  in de volgorde waarop ze bezocht worden.
- Als deel nummer  $i$  moet afgewerkt worden en je hebt een top met nummer  $m$  op diepte  $d$  bereikt dan ga door naar de volgende diepte als en slechts als  $i = m \pmod k$ . Anders backtrack.
- Als de delen  $0, 1, \dots, k - 1$  afgewerkt zijn (dat kan in parallel gebeuren) is de hele recursieboom doorzocht. Communicatie tussen de delen is alleen het mededelen van nieuwe bovengrenzen.

Deze manier om het doorzoeken van een recursieboom op te splitsen in deeltaken zien jullie in Figuur 19. Diepte 2 en maar 16 toppen op de diepte is

niet echt realistisch (het moeten er **veel** meer zijn), maar realistische gevallen kan je lang niet meer tekenen. . .

**Maar:** je moet erop letten dat in het gedeelte dat in elk deel wordt doorlopen – dus de niveaus ten hoogste de diepte waar gesplitst wordt – niet gesnoeid wordt of altijd met dezelfde waarde voor de grens gesnoeid wordt – anders is de nummering van de toppen op het splitlevel verschillend in delen die snel klaar zijn en delen die laat klaar zijn en misschien al met een betere grens kunnen werken.

Maar let op: het argument dat computers die gemakkelijke delen draaien ook sneller een nieuw deel opstarten geldt hier niet meer – wij moeten gewoon hopen dat door het feit dat elk deel onderdelen bevat ook de moeilijke gevallen (de takken met grote complexiteit) goed verdeeld zijn.



Figuur 19: De recursieboom voor het geval dat op diepte  $d = 2$  gesplitst wordt en deel nummer 2 voor  $k = 5$  afgewerkt wordt. Voor toppen met een nummer dat gelijk is aan  $2 \pmod{5}$  wordt doorgegaan en voor de anderen gestopt.

Inderdaad is het niet echt nodig dat je altijd op dezelfde diepte splitst. Soms is het beter in sommige takken op een grotere diepte te splitsen dan in andere. De bedoeling is gewoon dat je een verzameling van toppen in de boom hebt die een relatief klein gedeelte bij de wortel splitst van een groot gedeelte en waar je elke van deze splitstoppen vanuit het gedeelte dat altijd doorlopen wordt kan bereiken. Altijd op dezelfde diepte te splitsen is gewoon één van heel veel mogelijkheden die hier als voorbeeld werd gebruikt.

Over de tijd die nodig is om alle delen af te werken in vergelijking met de tijd als je het in één deel laat draaien, kan niet veel bewezen worden. Omdat de volgorde waarin de boom wordt doorlopen sterk verschilt van de volgorde waarin het in één deel wordt doorlopen – en dus de verschillende takken met andere waarden voor de bovengrens worden doorlopen, kan de som van de tijden theoretisch zelfs arbitrair veel kleiner zijn dan de tijd als je het in één deel laat draaien. Dat is geen strijdigheid met de observatie dat  $t^s(n)$  in feite

$O(p * t_p(n))$  is omdat  $t^s(n)$  betrekking heeft op een optimaal algoritme. – en in dit geval zou een optimaal algoritme de takken in een andere volgorde doorlopen.

Aan de andere kant kan het theoretisch ook gebeuren dat elk deel bijna even veel tijd vraagt als de hele taak. . .

In de praktijk zal deze manier van opsplitsen een goede manier zijn om dingen in het parallel op te starten waarbij de som van de tijden ongeveer even groot is als de tijd als het serieel wordt opgestart.

**Oefening 114** *Ons argument dat een serieel algoritme altijd ten hoogste  $p$  keer trager is dan een parallel algoritme was dat een serieel algoritme de parallele berekeningen kan simuleren. Beschrijf **expliciet** wat dat voor dit verdeeld algoritme betekent. Hoe werkt dit serieel algoritme voor een gegeven  $p$ ? Je mag veronderstellen dat  $p$  ook het aantal delen is – dus alle delen in het parallel worden opgestart.*

**Oefening 115** *Voor het volgende mag je stellen dat als een optimum gevonden wordt het als resultaat heeft dat alle nog lopende delen onmiddellijk kunnen stoppen omdat de snoeicriteria merken dat deze waarde niet verbeterd kan worden. Deze veronderstelling is in veel gevallen inderdaad realistisch (bv. als een langste cykel in een graaf gezocht wordt en een Hamiltoniaanse cykel wordt gevonden).*

*Stel bovendien dat alle delen tegelijk worden opgestart.*

- Gegeven een constante  $C \geq 1$ . Schets een scenario (precies: de structuur van een recursieboom) waar voor een branch and bound algoritme dat op de geziene manier in delen wordt opgesplitst de som van de tijden  $C$  keer kleiner is dan de tijd als het in één deel wordt opgestart.
- Gegeven een constante  $C \geq 1$ . Schets een scenario (precies: de structuur van een recursieboom) waar voor een branch and bound algoritme dat op de geziene manier in delen wordt opgesplitst de som van de tijden  $C$  keer groter is dan de tijd als het in één deel wordt opgestart.

**Oefening 116** *Herschrijf de pseudocode in Algoritme 17 zo dat het branch and bound algoritme op de net geziene manier opgesplitst wordt.*

**Oefening 117** *Gegeven gewichten  $g_1 \leq 1, \dots, g_n \leq 1$ . Gezocht is het minimale aantal vrachtwagens met capaciteit 1 dat nodig is om deze gewichten te vervoeren.*



- *Schrijf de pseudocode – met de nodige delen om het algoritme met de modulo optie te kunnen verdelen – voor een branch and bound algoritme dat de gewichten  $g_1, \dots, g_n$  in deze volgorde op elke mogelijke manier plaatst en de waarde  $b$  van een beste oplossing tot nu toe altijd bijhoudt.*

*Als bounding criterium gebruik het volgende: als de som  $k$  van alle gewichten die nog niet geplaatst zijn en op geen van de al gebruikte  $a$  vrachtwagens past, voldoet aan  $a + k > b - 1$ , kan je snoeien. Je mag veronderstellen dat een functie `get_k()` al bestaat die de waarde van  $k$  teruggeeft.*

- *Geef ook een andere manier splitstoppen te bepalen dan altijd op dezelfde diepte te splitsen. Houdt er bv. rekening mee hoeveel vrachtwagens al gebruikt zijn.*

**Oefening 118** *Stel dat je pas nadat al een groot deel van de taken afgewerkt is, merkt dat één deel toch veel langer draait dan alle anderen. Geef een oplossing hoe je deze taak verder kan opsplitsen zonder de andere delen die al afgewerkt zijn opnieuw op te starten.*

**Oefening 119** *Hetzelfde principe kan ook gebruikt worden voor het doorzoeken van recursiebomen waar geen optimum gezocht wordt, maar waar output op verschillende levels van de zoekboom gebeurt. Waarop moet je hier letten om nog altijd dezelfde verzameling van outputs te hebben en hoe los je het probleem op?*

## 6.2 Op weg met MPI

Terwijl de rest van het hoofdstuk over parallele algoritmen vrij theoretisch is en (nog) niet praktisch toepasbaar, zullen wij het hier over iets hebben dat wel heel belangrijk is voor toepassingen.

De bedoeling van dit hoofdstuk is niet, een echte inleiding te geven in MPI. De bedoeling is jullie op weg te helpen, MPI in C-programma's te gebruiken – de rest van de weg kunnen jullie dan gemakkelijk zelf vinden. MPI is ook geen algoritme, zodat een uitgebreide inleiding hier ook niet op de juiste plaats zou zijn. MPI (Message Passing Interface) is een bibliotheek van functies die jullie kunnen gebruiken om verschillende processen met elkaar te laten communiceren. MPI bestaat niet alleen voor C, maar onze voorbeelden hier zullen allemaal voor C zijn. Veel van de theoretische modellen die wij later zullen zien, kunnen *in principe* met MPI geïmplementeerd worden – maar in de meeste gevallen niet met de in de theoretische modellen

veronderstelde complexiteit. Het grote voordeel van MPI is, dat het een standaard is. De MPI-1 standaard werd al in 1994 vastgelegd en de bibliotheek staat ondertussen voor alle belangrijke platformen ter beschikking. Een met MPI geparalleliseerd algoritme is dus portable als je buiten de MPI-functies geen syteemspecifieke dingen doet: jullie kunnen het algoritme op een laptop ontwikkelen en zonder wijzigingen bv. op een linux-cluster laten draaien of op een machine met een groot aantal cores. Het is de taak van de MPI-ontwikkelaars de topologie van het netwerk optimaal te gebruiken. Ook al zijn de programma's portable – de manier waarop de programma's **opgestart** moeten worden, kan verschillen. Die hangt er ook van af, of bv. rekening moet gehouden worden met vrije processoren, queues, etc. In deze lesnota's wordt de manier gegeven die je op één linux computer kan gebruiken (nadat je MPI geïnstalleerd hebt).

### 6.2.1 MPI opstarten

Maar het beste is als wij onmiddellijk beginnen. Kijk naar het volgende eenvoudige programma:

#### Programma 1

```
// programma simpleworld.c
#include <mpi.h>
#include <stdio.h>
#include <unistd.h>

int main (int argc, char* argv[])
{
    MPI_Init(&argc, &argv);      /* starts MPI */
    sleep(10);                   /* body */
    printf( "Hello world!\n");    /* body */
    MPI_Finalize();              /* closes MPI */
}
```

Belangrijk is hier de `#include <mpi.h>` in de header die ervoor zorgt dat alle MPI-functies gekend en goed gedefinieerd zijn. Elk MPI-programma start bovendien met `MPI_Init(&argc, &argv)` en stopt met `MPI_Finalize()`. Dit programma kan je nu met `mpicc -o simpleworld simpleworld.c` compileren. Als jullie dat nu parallel met 4 processen willen laten draaien, dan kan dat met `mpirun -n 4 simpleworld`. Daarbij staat `-n x` ervoor dat je `x` processen wilt opstarten. Deze processen worden afhankelijk van waar

je werkt, op dezelfde computer opgestart of op verschillende andere computers. Als je het op een alleenstaande Linux-computer opstart, dan geeft de `sleep(10)` je voldoende tijd om te zien dat er daadwerkelijk 4 processen `simpleworld` aan het draaien zijn. Na 10 seconden zie je dan het resultaat: 4 lijnen met `Hello world!`. In oude implementaties van MPI kan het gebeuren dat niet alle processen `stdout` en `stderr` kunnen gebruiken, maar “meestal” zullen `stdout` en `stderr` naar de `stdout` en `stderr` van het proces `mpirun` op de masternode omgeleid worden. Ook al is de kans klein dat jullie ooit een implementatie tegenkomen die dat niet doet: zelfs in de MPI 3 standaard is niet vastgelegd dat elke geldige implementatie dat **moet** doen. Wie gegarandeerd volledig portable programma’s wil schrijven, stuurt best alle data naar het root proces (welk proces dat is, zullen wij nog zien) en laat dat de output doen.

Bijna alle MPI-functies geven een (int-) errorwaarde terug – ook `MPI_init` en `MPI_Finalize`. Hier willen wij het programma zo eenvoudig mogelijk houden, maar een goed programma zou deze waarde natuurlijk het best toetsen, bv met

```
resultaat = MPI_Init (&argc, &argv);  
if (resultaat != MPI_SUCCESS).....
```

waarbij `MPI_SUCCESS` een MPI-constante is die aangeeft dat er geen fout opgetreden is. Het testen van errorcodes is belangrijk voor een vaak gebruikt programma, maar hier zou het gewoon afleiden van de hoofdzaak. De programma’s hier zijn ten slotte alleen maar bedoeld om een beetje te verstaan wat er gebeurt.

Dat is al een begin, maar natuurlijk niet echt boeiend. Er is bv. nog helemaal geen communicatie en alle processen doen hetzelfde. Om ervoor te zorgen dat verschillende processen verschillende dingen kunnen doen, moeten de processen een soort id hebben op basis waarvan ze kunnen beslissen wat ze moeten doen. Een eerste aanzet zou kunnen zijn met `getpid()` de proces id te gebruiken, maar als de processen op verschillende machines draaien, zouden die zelfs voor verschillende processen identiek kunnen zijn en bovendien zouden we niet op voorhand weten, welke nummers toegewezen worden. Dat deugt dus niet.

MPI gebruikt het concept van een communicator, die een groep van processen voorstelt die met elkaar kunnen communiceren. Een communicator is als datatype door middel van `typedef` gedefinieerd. Het datatype heet `MPI_Comm` en een belangrijke **constante** met dit datatype is `MPI_COMM_WORLD` – de constante groep die alle processen bevat die in het begin worden opgestart. Inderdaad kan het aantal processen op een dynamische manier groeien. Als

dat gebeurt, beschrijft de constante `MPI_COMM_WORLD` niet meer alle processen en kunnen de verschillende nieuwe processen zelfs disjuncte verzamelingen als *world* beschrijven. In deze mini-inleiding zal het aantal processen altijd constant blijven.

Belangrijke functies voor een dergelijke communicator die jullie zeker nodig zullen hebben, zijn:

```
MPI_Comm_rank (MPI_Comm communicator, int *index);
/* bepaal de index van dit proces in "communicator"*/

MPI_Comm_size (MPI_Comm communicator, int *grootte);
/* bepaal het aantal processen in "communicator"*/
```

Ook deze functies geven een errorcode terug. Daarom moeten de adressen van de integer variabelen waar de opgevraagde waarden naartoe geschreven moeten worden, meegegeven worden. De index in een communicator is een goede id voor het proces die we kunnen gebruiken om te beslissen wat een proces doet. Als er  $k$  processen in een communicator zijn, dan zijn de indexen  $0, \dots, k-1$ . De index 0 in `MPI_COMM_WORLD` is een bijzonder proces dat ook het root proces wordt genoemd. In `Open_MPI` en andere implementaties volgens `MPI 2` standaard wordt de `stdin` van `mpirun` aan dit proces doorgegeven en alle andere processen hebben `/dev/null` als `stdin`. Wij zullen data altijd van `stdin` lezen (bv. met `fread()` of `scanf()`) maar voor de duidelijkheid in het programma er expliciet voor zorgen dat alleen proces 0 probeert te lezen.

Met deze functies kunnen wij ons programma al een beetje interessanter maken:

## Programma 2

```
#include <mpi.h>
#include <stdio.h>

void zeg_niets()
{ printf("Ik zeg niet wie ik ben!\n"); }

void zeg_iets(int mijn_id, int totaal)
{ printf("Ik ben proces nummer %d van totaal %d !\n",mijn_id,totaal); }

int main (int argc, char* argv[])
{
```

```

int mijn_id, grootte, i;

MPI_Init (&argc, &argv);                /* starts MPI */
MPI_Comm_rank (MPI_COMM_WORLD, &mijn_id); /* vraag mijn nummer op */
MPI_Comm_size (MPI_COMM_WORLD, &grootte); /* hoeveel processen zijn er? */
if (mijn_id%2) zeg_niets();
else zeg_iets(mijn_id,grootte);
MPI_Finalize();
}

```

Als wij dat compileren en met `mpirun -n 6` laten draaien, **kan** de output bv als volgt zijn:

```

Ik ben proces nummer 0 van totaal 6 !
Ik zeg niet wie ik ben!
Ik zeg niet wie ik ben!
Ik ben proces nummer 4 van totaal 6 !
Ik zeg niet wie ik ben!
Ik ben proces nummer 2 van totaal 6 !

```

Alle processen draaien tegelijk en welk proces eerst zijn resultaat doorgeeft, is niet vastgelegd. De verschillende MPI-processen hebben geen synchronisatie (als die niet door zekere communicaties wordt opgelegd) en draaien onafhankelijk van elkaar. Inderdaad kan de output van een proces zelfs onderbroken zijn door andere output. Een goed programma zou er dus het best zelf voor zorgen dat de output door één proces gesynchroniseerd wordt. Maar dat kunnen wij op dit ogenblik nog niet. Wij kunnen wel al verschillende processen verschillende dingen laten doen – gebaseerd op hun index in een communicator.

Maar toch al is er nog geen communicatie...

**Oefening 120** *Schrijf een MPI-programma waar elk proces een lijst van identificaties van de andere processen uitvoert met wie het via de `MPI_COMM_WORLD` communicator kan communiceren.*

## 6.2.2 Berichten sturen en ontvangen

De meest eenvoudige vorm van communicatie is de communicatie tussen twee vastgelegde processen (vaak point-to-point communicatie genoemd). Daarover zullen wij het dus eerst hebben.

Een bericht heeft twee belangrijke delen: de enveloppe en de inhoud. De enveloppe bevat informatie over waar het bericht naartoe gaat, hoe groot die

MPI datatype	C datatype
MPI_CHAR	char
MPI_SHORT	short int
MPI_INT	int
MPI_LONG	long int
MPI_LONG_LONG_INT	long long int
MPI_SIGNED_CHAR	char
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED_SHORT	unsigned short
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long int
MPI_UNSIGNED_LONG_LONG	unsigned long long int
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double

Tabel 1: Sommige elementaire datatypes en de bijbehorende C datatypes.

is, etc. en de inhoud is dat wat je echt wilt sturen. In principe is de inhoud altijd een array, maar MPI geeft je ook de mogelijkheid afgeleide datatypes te definiëren – net zoals structs in C. Wij zullen alleen met de elementaire datatypes werken, waarvan we de belangrijkste in tabel 1 geven.

Een typische inhoud van een bericht is dus bv. een array van ints. Belangrijk is daarbij dat de inhoud tijdens het sturen soms *vertaald* moet worden. Als je bv. op een cluster werkt waar zowel Little-Endian als Big-Endian machines zijn en je stuurt een int, dan wil je dat de machine die het bericht ontvangt hetzelfde getal ontvangt – MPI zal de volgorde zo wijzigen dat de int op de andere machine hetzelfde getal voorstelt. Moeilijker wordt het als de datatypes verschillende groottes hebben – dan moet MPI soms afronden. Dat zijn allemaal problemen die MPI voor je oplost – omdat je niet op voorhand weet op welk soort machine welk proces draait, **kan** je de conversies niet zelf doen.

Een manier voor het sturen van point-to-point communicatie is de functie

```
int MPI_Send(void *buf, int aantal, MPI_Datatype type,
             int ontvanger, int kenmerk, MPI_Comm comm)
```

De inhoud van ons bericht is beschreven door de 3 variabelen `buf`, `aantal` en `type`, waarbij je `aantal`, en `type` als deel van de enveloppe kan rekenen, omdat het de inhoud beschrijft maar niet de inhoud is. In deze functie is

**\*buf** een pointer naar de array die de inhoud bevat die we willen sturen, **aantal** het aantal elementen en **type** het datatype van de elementen van de array.

De variabele **ontvanger** is het nummer van het proces in de communicator **comm** die het bericht moet ontvangen. Als laatste deel van de enveloppe heb je nog de variabele **kenmerk** die je als een identificatie van het bericht kan zien. Als een proces  $p_0$  bv. meerdere berichten aan hetzelfde proces  $p_1$  stuurt, dan kunnen die door de kenmerken nog herkend worden en de ontvanger kan precies die berichten ophalen die hij op dat moment nodig heeft.

Aan de andere kant moet het bericht ook ontvangen worden. Daarvoor heb je de functie

```
int MPI_Recv(void *buf, int aantal, MPI_Datatype type, int zender,
             int kenmerk, MPI_Comm comm, MPI_Status *status)
```

Als er een bericht is van het proces met nummer **zender** in de communicator **comm** met identificatie **kenmerk**, dan worden **aantal** items van type **type** naar **buf** geschreven. Het is **jouw** taak om het programma zo te schrijven dat **aantal** en **type** overeenkomen met de variabelen in de functie **MPI\_Send()** die het bericht met deze identificatie heeft gestuurd. Als dat niet zo is, gebeurt er een fout, die soms door MPI herkend kan worden (bv. *truncated message* – te weinig bytes gelezen) maar niet altijd.

In de struct status heb je altijd **MPI\_SOURCE** – het nummer van het proces dat het bericht heeft gestuurd, **MPI\_TAG** – de identificatie van het bericht en **MPI\_ERROR** – de error status. Op het eerste gezicht lijkt een dergelijke struct overbodig omdat het alleen maar informatie bevat die we al kennen, maar je kan ook wildcards gebruiken voor **zender** (**MPI\_ANY\_SOURCE**) en **kenmerk** (**MPI\_ANY\_TAG**) en dan kan je met de status-struct deze informatie uitvissen. Een wildcard **MPI\_ANY\_DEST** bestaat **niet**. Als je een bericht stuurt, moet dus duidelijk zijn waar dat naartoe gaat.

De werking kan bv. in het volgende programma gezien worden:

### Programma 3

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>

void proc(int id)
{
    int number, i;
    MPI_Status status;
```

```

    if (id==1)
    { i=0;
      do
      {
        MPI_Recv((void *)&number, 1, MPI_INT,0, i, MPI_COMM_WORLD,&status);
        fprintf(stderr,"Process %d received number %d with tag %d \n",
                  id,number,status.MPI_TAG);
        i= (i+2)%5;
      }
      while (i!=0);
    }
    else fprintf(stderr,"Process %d has nothing to do.\n",id);
    return;
}

void proc0(int id)
{
    int i,j;

    for (i=0;i<5;i++)
    { j=3*i;
      MPI_Send((void *)&j, 1, MPI_INT,1, i, MPI_COMM_WORLD);
      fprintf(stderr,"Process %d sent number %d\n",id,j);
    }
    return;
}

int main (int argc, char* argv[])
{
    int myrank;

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &myrank);

    if (myrank==0) proc0(myrank);
    else proc(myrank);

    MPI_Finalize();
    return 0;
}

```



De rare vorm waarop de do-while lus in dit voorbeeld wordt doorlopen, is om te zien dat de berichten inderdaad niet op een FIFO of LIFO manier aan een queue voor een proces worden toegevoegd, maar op basis van de identificatie opgehaald kunnen worden. Als je dit programma met 3 processen opstart, is de output dus bv.:

```
Process 2 has nothing to do.  
Process 0 sent number 0  
Process 0 sent number 3  
Process 0 sent number 6  
Process 0 sent number 9  
Process 0 sent number 12  
Process 1 received number 0 with tag 0  
Process 1 received number 6 with tag 2  
Process 1 received number 12 with tag 4  
Process 1 received number 3 with tag 1  
Process 1 received number 9 with tag 3  
Process 3 has nothing to do.
```

De volgorde van de output van de verschillende processen kan verschillend zijn (bv. `Process 2 has nothing to do.` hoeft niet het eerste bericht te zijn, maar kan ook later komen), maar voor hetzelfde proces zal de volgorde van berichten altijd dezelfde zijn. Een proces kan dus meerdere berichten hebben die wachten om opgehaald te worden.

**Oefening 121** *Schrijf een programma voor 3 processen dat als volgt werkt: De processen 1 en 2 genereren toevallige getallen met `random()`, dat geïnitialiseerd wordt met `srandom(time(0)+process_id)`. Reeksen van oneven getallen worden in een array bijgehouden en zodra een even getal opduikt, geldt de reeks als afgesloten. Als de reeks niet leeg is, wordt hij naar proces 0 gestuurd. Als de reeks van getallen dus 4,1,6,8,5,7,9,8,11... is, wordt dus eerst een array met alleen maar 1 doorgestuurd, dan een array met 5,7,9, etc.*

*Proces 0 leest afwisselend de berichten van proces 1 en 2 en houdt voor beide processen 1 en 2 bij wat de som van de lengten van de doorgestuurde reeksen is en wat de langste reeks was. Zodra één van de processen in totaal 1.000.000 getallen heeft doorgestuurd, schrijft proces 1 de langste reeks van dit proces uit en het programma stopt.*

*Let op: proces 0 weet op voorhand niet hoe lang de reeks van oneven getallen is, je moet er dus voor zorgen dat hij weet hoeveel hij moet lezen. Als proces 0 merkt dat het resultaat gevonden is, gaan de andere twee processen nog altijd door. Zorg ervoor dat die ook stoppen.*

Beide functies – `MPI_Send()` en `MPI_Recv()` moeten eerst *afgewerkt* zijn voordat het programma dat ze gebruikt kan doorgaan – ze blokkeren dus de voortgang van het programma en kunnen tot een deadlock leiden. Het is dus belangrijk precies te weten wat *afgewerkt* in dit geval betekent, dus **wanneer** een oproep van deze functies *afgewerkt* is. Inderdaad kan ook Programma 3 in principe in een deadlock terechtkomen – zie Oefening 122.

Wat gebeurt er bv. als je het volgende programma met 3 processen opstart?

#### Programma 4

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char* argv[])
{
    int mijn_id, grootte, i, message[1000], message2[1000];
    MPI_Status status;

    MPI_Init (&argc, &argv);          /* starts MPI */
    MPI_Comm_rank (MPI_COMM_WORLD, &mijn_id);
    MPI_Comm_size (MPI_COMM_WORLD, &grootte);

    for (i=0;i<1000;i++) message[i]=mijn_id;

    MPI_Send((void *)message, 1000, MPI_INT, (mijn_id+1)%grootte,
             1,MPI_COMM_WORLD);
    MPI_Recv((void *)message2, 1000, MPI_INT,MPI_ANY_SOURCE,
             MPI_ANY_TAG, MPI_COMM_WORLD,&status);

    fprintf(stderr,"Process %d received list of numbers %d with tag %d \n",
            mijn_id,message2[0],status.MPI_TAG);

    MPI_Finalize();
    return 0;
}
```

Het resultaat op mijn computer (met 3 processen) was:

```
Process 0 received list of numbers 2 with tag 1
Process 1 received list of numbers 0 with tag 1
Process 2 received list of numbers 1 with tag 1
```

Alles lijkt dus te draaien zo als gehoopt, maar op een andere computer kan het gebeuren dat het programma in een deadlock raakt. Zelfs op een computer waar het lijkt te werken, kan het gebeuren dat het niet meer draait als je meer data stuurt dan maar 1000 integers. De reden daarvoor is de manier waarop het versturen werkt. Voor het ontvangen (`MPI_Recv()`) is het duidelijk wat er gebeurt: het programma wacht totdat het opgevraagde bericht er is, leest het gegeven aantal items (als dat kan – anders is er een fout) en dan is de functie afgewerkt en het programma kan doorgaan.

Voor het sturen (`MPI_Send()`) is dat niet zo duidelijk. Wanneer is de functie klaar – dus wanneer is een bericht *gestuurd*? MPI kan (en zal) dat op meerdere manieren doen. MPI kan het bericht in een eigen buffer schrijven, wat heel snel kan gebeuren. Zodra MPI het bericht in zijn eigen buffer heeft, is `MPI_Send()` gedaan en kan het programma dat iets wil sturen doorgaan. Maar als MPI niet voldoende eigen buffer heeft, kan het ook wachten totdat het proces dat het bericht wil ontvangen dat afhaalt en dan de inhoud van de array zonder bufferen van het ene proces naar het andere schrijven. Dan kan het dus gebeuren dat een proces heel lang op een andere computer moet wachten – en zelfs erger: je kan een deadlock krijgen. Als niet gebufferd wordt – als designbeslissing van de speciale MPI-implementatie of omdat er te veel data gestuurd wordt – zullen in Programma 4 alle processen tegelijk beginnen te schrijven en zolang de andere processen het bericht niet ophalen, zal geen van de processen doorgaan. Maar omdat **alle** processen wachten, zal geen proces iets ophalen en ontstaat er een deadlock...

**Oefening 122** *Beschrijf waar en hoe Programma 3 in principe in een deadlock terecht kan komen.*

**Oefening 123** *Schrijf een MPI-programma dat dezelfde berichten aan dezelfde processen stuurt als Programma 4, maar waar er verzekerd is dat geen deadlock kan ontstaan.*

*Bespreek ook (informeel) de tijd die Programma 4 (als er geen deadlock is) en jouw programma nodig hebben (zowel de tijd per proces als ook de totale tijd vanaf het opstarten totdat alle processen klaar zijn).*

**Oefening 124** *Implementeer een – wel niet erg snelle – manier om  $\sum_{i=1}^n i$  te berekenen:*

*Je werkt met  $n + 1$  processen. Proces  $n$  stuurt zijn id aan proces  $n - 1$ . De processen  $k \in \{1, \dots, n - 1\}$  ontvangen een getal van proces  $k + 1$ , tellen hun id erbij op en sturen het resultaat naar proces  $k - 1$ . Ten slotte ontvangt proces 0 een getal van proces 1 en schrijft die als resultaat.*

**Oefening 125** *Schrijf een MPI-programma voor mergesort. Proces 0 leest een reeks van getallen van `stdin` in en slaat die in een array op. De getallen staan binair in een bestand dat je door middel van een pipe als `stdin` van `mpirun` gebruikt. Als je  $k$  processen hebt, wordt de array eerst in  $k$  delen van ongeveer gelijke grootte opgedeeld en worden die aan de verschillende processen doorgestuurd (proces 0 houdt natuurlijk één). De verschillende processen sorteren de delen en sturen die door naar andere processen om te mergen, zodat elk proces ten hoogste 2 delen heeft om te mergen. Als ten slotte nog maar 1 proces bezig is (zorg ervoor dat dat proces 0 is), merget die zijn twee delen en schrijft het resultaat naar `stdout`.*

### 6.2.3 Berichten sturen en ontvangen – zonder de voortgang te blokkeren

Ook als de communicatie goed georganiseerd is, zit je bij de methode zoals wij die tot nu toe gezien hebben met een probleem: soms weet je gewoon niet of er een bericht voor je is! Een typisch geval zijn onze verdeelde branch and bound algoritmen: daar weet een proces dat er een betere grens van een ander proces *zou kunnen zijn* – maar misschien is die er ook niet, en als die er niet is, kan de ontvanger zonder problemen doorgaan. Als je voor een dergelijke toepassing `MPI_Recv()` zou gebruiken, zou je niet alleen tijd verspillen als er geen bericht klaarligt – het zou zelfs kunnen dat het proces nooit stopt omdat er geen bericht komt.

Voor dergelijke toepassingen heb je de functies `MPI_Isend()` en `MPI_Irecv()` waarbij de “I” voor *initiate* staat. Beide functies starten een communicatieoperatie en keren terug nog voordat die afgewerkt is. Om later naar deze opdracht te kunnen refereren, geven ze een MPI handle terug. De manier waarop de functies opgeroepen worden is als volgt:

```
int MPI_Isend(void *buf, int aantal, MPI_Datatype type,
              int ontvanger, int kenmerk, MPI_Comm comm,
              MPI_REQUEST *opdrachtnummer)

int MPI_Irecv(void *buf, int aantal, MPI_Datatype type, int zender,
              int kenmerk, MPI_Comm comm,
              MPI_REQUEST *opdrachtnummer)
```

De variabele `opdrachtnummer` is hier de plaats waar de MPI handle wordt naartoe geschreven. De functies lijken – op de MPI handle na – sterk op de blokkerende functies. In `MPI_Irecv()` ontbreekt wel de variabele `status`, maar die kan natuurlijk pas ingevuld worden als de ontvangsoperatie afgewerkt is.

Of een operatie klaar is, kan je testen met de functie `MPI_Test()` die de volgende parameters nodig heeft:

```
int MPI_Test(MPI_Request *opdracht, int *flag, MPI_Status *status)
```

Als `flag true` is, is de opdracht met handle *\*opdracht* klaar – anders niet. Als de opdracht klaar is, staat in `status` wat je ook na een `MPI_Recv()` in `status` zou vinden . Na een `MPI_Isend()`-opdracht **kan** in `status` de errorcode van de `MPI_Isend()`-operatie staan.

De werking kan in het volgende programma gezien worden:

### Programma 5

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>

void proc(int id)
{
    int number, gedaan, seconden=0;
    MPI_Status status;
    MPI_Request opdracht;

    MPI_Irecv((void *)&number, 1, MPI_INT, MPI_ANY_SOURCE, MPI_ANY_TAG,
              MPI_COMM_WORLD, &opdracht);

    for (gedaan=0; !gedaan; )
    {
        MPI_Test(&opdracht, &gedaan, &status);
        if (gedaan)
            fprintf(stderr, "Proces %d heeft bericht van %d ontvangen \n",
                    id, status.MPI_SOURCE);
        else { sleep(5); seconden+=5;
              fprintf(stderr, "Proces %d heeft  %d seconden gewacht!\n", id, seconden);
            }
    }
    return;
}

void proc0(int id, int total)
{
    int number=1, gedaan;
```

```

MPI_Status status;
MPI_Request opdracht;

if (total>1)
{ sleep(13);
  MPI_Isend((void *)&number, 1, MPI_INT,1, 0, MPI_COMM_WORLD,&opdracht);
  MPI_Test(&opdracht,&gedaan,&status);
  if (gedaan) fprintf(stderr,"Proces %d heeft nummer %d onmiddelijk gestuurd!\n",
                      id,number);
}
return;
}

int main (int argc, char* argv[])
{
  int myrank, size;

  MPI_Init (&argc, &argv);
  MPI_Comm_rank (MPI_COMM_WORLD, &myrank);
  MPI_Comm_size (MPI_COMM_WORLD, &size);

  if (myrank==0) proc0(myrank,size);
  else proc(myrank);

  MPI_Finalize();
  return 0;
}

```

Opgestart met 3 processen zou de output bv. als volgt kunnen zijn:

```

Proces 1 heeft  5 seconden gewacht!
Proces 2 heeft  5 seconden gewacht!
Proces 1 heeft 10 seconden gewacht!
Proces 2 heeft 10 seconden gewacht!
Proces 0 heeft nummer 1 onmiddelijk gestuurd!
Proces 1 heeft 15 seconden gewacht!
Proces 2 heeft 15 seconden gewacht!
Proces 1 heeft bericht van 0 ontvangen
Proces 2 heeft 20 seconden gewacht!
Proces 2 heeft 25 seconden gewacht!
.

```

.  
.

Proces 1 en 2 initialiseren een ontvangstopdracht en gaan – zoals de output toont – door met hun werk als het bericht er nog niet is. In dit geval zou proces 2 oneindig doorgaan omdat er natuurlijk geen bericht komt...

**Belangrijk** is dat een proces dat `MPI_Isend()` gebruikt de variabele waarin de inhoud is opgeslaan niet wijzigt voordat gedetecteerd is dat het versturen klaar is. Pas als `MPI_Test()` zegt dat het proces klaar is mag die gewijzigd worden – anders kan een fout optreden.

De `MPI_I...()` routines hebben duidelijk voordelen tegenover de blokkerende routines: Ze kunnen alles wat de blokkerende routines kunnen en meer. Bovendien heb je gegarandeerd geen deadlock. Aan de andere kant zijn programma's met de niet blokkerende routines moeilijker te lezen en te verstaan, omdat het effect van het sturen en ontvangen niet zo duidelijk lokaal afgebakken is – bv. dezelfde ontvangstopdracht kan in verschillende functies getoetst moeten worden.

#### 6.2.4 Verder met MPI?

Dit waren slechts de beginselen van MPI, maar met deze weinige basisoperaties kan je al zinvolle en nuttige programma's schrijven. De gehele functionaliteit van MPI is natuurlijk veel groter, maar nu jullie op weg zijn, kunnen jullie zeker zelf verder gaan. Andere nuttige functies (die je wel met de al geziene functies ook zelf kan implementeren) zijn bv. `MPI_Bcast()` (*broadcasting*), `MPI_Gather()` (*gathering*) of `MPI_Scatter()` (*scattering*).

### 6.3 De modellen voor parallel computing

In dit deel zullen wij nu sommige theoretische modellen voor computers zien die – dat moet wel toegegeven worden – op dit moment (nog?) niet als hardware verwezenlijkt zijn. Maar zoals al in het begin gezegd, is niet te voorspellen of de multicore processoren de parallelle algoritmen niet nieuw leven inblazen en net dit deel voor jullie in toekomst bijzonder nuttig zal zijn... Zeker is dat er veel nieuwe ideeën voor parallelle algoritmen ontwikkeld werden. En ideeën zijn natuurlijk altijd interessant!

De modellen waarover wij het hier hebben, zijn wiskundige modellen voor parallelle computers. Terwijl er voor de gewone computer een door iedereen aanvaard model bestaat – het RAM-model – bestaan voor parallelle computers meerdere verschillende modellen en sommige algoritmen zijn alleen

bedoeld voor één model en werken niet op andere modellen zonder belangrijke wijzigingen.

Wij bespreken alleen maar heel kort modellen waar de processoren deel uitmaken van een netwerk dat beschrijft welke processoren met elkaar kunnen communiceren en welke dat niet kunnen. Dergelijke modellen waren een tijdje geleden heel populair en in de vorm van transputers bestonden die zelfs als hardware. Transputers konden zelfs zo geconfigureerd worden dat je een netwerk hebt dat speciaal voor het algoritme dat erop moet draaien geschikt is. Maar er was ook veel onderzoek over netwerken met goede algemene eigenschappen – bv. de *hypercubes* – en over het simuleren van algemene netwerken op specifieke netwerken.

Wij zullen het hier vooral over variaties van het eenvoudigste – en misschien belangrijkste – model hebben. Het is de meest directe veralgemening van het RAM-model en heet **Parallel Random Access Machine** – PRAM.

In dit model heb je gewone processoren die je bv. door een nummer kan identificeren – dus bv. processoren  $p_1, \dots, p_k$ . Elke processor heeft zijn eigen programma (dat heet *Multiple Instruction*). Sommige algoritmen gebruiken ook het *Single Instruction* model waar alle processoren hetzelfde programma moeten gebruiken, maar wij veronderstellen gewoon het algemenere concept van Multiple Instruction waar de programma's van de processoren gelijk mogen zijn maar het niet vereist is. De processoren zijn bovendien gesynchroniseerd – dat betekent dat je op verschillende punten van een programma mag veronderstellen dat de processoren van dit punt tegelijk vertrekken. Als je bv. een lus in de pseudocode hebt

```
for i=1 to 20 pardo { do_iets(i) }
```

(met **pardo** voor “*do in parallel*”) dan betekent dat dat voor  $1 \leq i \leq 20$  processor  $p_i$  de taak **do\_iets(i)** afwerkt en pas als alle processoren klaar zijn doorgegaan wordt – ook als de taak voor verschillende  $i$  langer duurt dan voor anderen. Als verschillende processoren dezelfde lijnen in een dergelijke lus afwerken wordt normaal ook verondersteld dat ze op hetzelfde moment met de lijnen beginnen – dat er dus ook synchronisatie binnen in de lus is. In ons geval zullen dat altijd maar sommige lijnen zijn en het zal duidelijk zijn waar de synchronisatie moet gebeuren.

Als in de pseudocode zoiets staat als

```
for i=1 to log n pardo { do_iets(i) }
```

dan betekent dat niet dat er een centrale processor is die de andere processoren bestuurt maar dat de individuele processoren – die hun eigen nummer kennen – alleen dan iets doen als hun nummer ten hoogste  $\log n$  is. De volgende lijn is dus bv. deel van de individuele programma's van de processoren



– waarbij  $i$  het nummer van de processor is en wij ervoor moeten zorgen dat de processor  $n$  ook kent:

```
if (i <= log n) { do_iets(i) }
```

De processoren hebben een gemeenschappelijk geheugen dat ze allemaal kunnen gebruiken om in constante tijd te lezen en te schrijven. Je kan ook veronderstellen dat elke processor bovendien nog lokaal geheugen heeft maar dat komt natuurlijk op hetzelfde neer als globaal geheugen waarvan een deel alleen door de ene processor gebruikt wordt. Maar als meerdere processoren hetzelfde geheugen gebruiken is er natuurlijk een probleem: wat als meerdere processoren **tegelijk** iets willen lezen of schrijven? Daarvoor bestaan verschillende modellen. Eén model verbiedt gelijktijdig te lezen/schrijven – dat heet Exclusive Read resp. Exclusive Write en wordt met ER resp. EW afgekort. Een ander model laat het toe gelijktijdig te lezen/schrijven – dat heet dan Concurrent Read resp. Concurrent Write en wordt CR resp. CW afgekort.

Maar voor CW is er nog een probleem: wat gebeurt als meerdere processoren tegelijk *verschillende* waarden in een variabele willen schrijven? Wij leggen vast dat het dan toevallig is welke waarde in de variabele terecht komt. (Maar er bestaan ook modellen die dat anders vastleggen...) Het algoritme moet dus correct werken om het even welke processor erin slaagt zijn waarde te schrijven. In het geval van ER of EW is het de verantwoordelijkheid van het algoritme ervoor te zorgen dat nooit verschillende processoren dezelfde variabele tegelijk willen lezen of er naartoe willen schrijven. Een algoritme dat dat niet garandeert is gewoon ongeldig!

Wij hebben dus 4 verschillende modellen van PRAMs: EREW-PRAM, ERCW-PRAM, CREW-PRAM en CRCW-PRAM.

Maar in alle modellen moet gegarandeerd zijn dat in dezelfde stap nooit één processor een variabele leest en één processor in dezelfde variabele schrijft. Het probleem kunnen jullie in Algoritme 18 zien.

### Algoritme 18

```
x=0; a=0;  
for i=1 to 2 pardo  
  if (i=1) x=1; else a=x;
```

De stappen  $x=1$  en  $a=x$  worden tegelijk uitgevoerd. De vraag is dus of in deze *stap* op het moment dat  $x$  wordt gelezen de waarde al gewijzigd is of niet... Natuurlijk kan je de modellen ook aanpassen door toe te laten dat gelijktijdig gelezen en geschreven *mag* worden maar te eisen dat het resultaat van het algoritme daarvan niet *mag* afhangen. In onze algoritmen zullen wij

er gewoon altijd op letten dat nooit tegelijk (dat is: in dezelfde *stap* – waarbij *stap* natuurlijk niet echt precies gedefinieerd is) uit een variabele gelezen en in dezelfde variabele geschreven wordt.

Maar misschien is het het beste nu eens naar voorbeeldalgoritmen te kijken. Het eerste voorbeeld kopieert een waarde *a* zodat deze waarde achteraf *n* keer in een array *copy*[] staat. In feite kan de waarde zelfs in meer vakken staan – precies  $2^{\lceil \log n \rceil}$ . Dit aantal wordt ook in de array *aantal*[] teruggegeven – in  $2^{\lceil \log n \rceil}$  kopieën. Het algoritme is geschikt voor een EREW-PRAM en draait in tijd  $O(\log n)$ . Wij schrijven  $\lceil n/2 \rceil = 2^{\lceil \log(n/2) \rceil}$ .

**Algoritme 19** *Waarden kopiëren in een EREW-PRAM*

```
copywaarden(a,n,copy[],aantal[])

{
  for j=1 to  $\lceil n/2 \rceil$  pardo  n_array[j]= 0;
    // een array om n voor alle n processoren tegelijk leesbaar te maken

  copy[1]=a;
  n_array[1]=n;
  aantal[1]=1; // een hulpparray dat zegt hoeveel kopie\en er al zijn

  for (i=1; i<n; i=2*i) // deze lijn is heel informeel. Voor de precieze
    // vorm kijk naar het programma per processor
  { for j=1 to i pardo
    { copy[j+aantal[j]]=copy[j];
      n_array[j+aantal[j]]=n_array[j];
      aantal[j]=aantal[j+aantal[j]]=2*aantal[j];
    }
  }
  // na deze lus staat a ten minste in copy[1]...copy[n]
}
```

Deze keer zullen wij het voor de duidelijkheid nog eens anders schrijven. Het volgende Algoritme 20 geeft je het programma voor processor *i*. Het wordt opgestart voor  $2^{\lceil \log(n/2) \rceil}$  processoren en elke processor kent zijn eigen nummer “nummer”.

**Algoritme 20** *Waarden kopiëren in een EREW-PRAM – programma voor processor nummer*

```

copywaarden(a,n,copy[],aantal[])

{
  n_array[nummer]= 0;

  if (nummer=1)
  { copy[1]=a;
    n_array[1]=n;
    aantal[1]=1;
  }

  for (i[nummer]=1; (n_array[nummer]= 0) || (i[nummer]<n_array[nummer]);
                                             i[nummer]=2*i[nummer])
  // aan elke van de volgende binnenste lussen wordt
  // tegelijk begonnen -- wij hebben dus een globale synchronisatie
  { if (nummer <= i[nummer])
    { copy[nummer+aantal[nummer]]=copy[nummer];
      n_array[nummer+aantal[nummer]]=n_array[nummer];
      aantal[nummer]=aantal[nummer+aantal[nummer]]=2*aantal[nummer];
    }
  }
  // na deze lus staat a ten minste in copy[1]...copy[n]
}

```

Als  $n$  kopieën gemaakt moeten worden, vraagt dit algoritme tijd  $O(\log n)$  en  $2^{\lceil \log(n/2) \rceil}$  processoren – dus  $O(n)$  processoren.

Belangrijk is dus te zien dat ook al kunnen EREW machines niet alle waarden tegelijk lezen, je die wel in tijd  $O(\log n)$  kan verdelen. Dus kan een ER machine een CR machine met een verliesfactor van ten hoogste  $O(\log n)$  (waarbij  $n$  het aantal processoren is) simuleren door waarden te verdelen. Maar dat is niet zo triviaal als het misschien lijkt: in één stap kunnen  $n$  processoren  $n$  verschillende variabelen wijzigen. Als die waarden dan allemaal verdeeld moeten worden, heb je extra processoren nodig om dat in tijd  $O(\log n)$  te doen. Bovendien moet je ook garanderen dat geen twee van de extra processoren dezelfde variabele willen wijzigen...

Een serieel algoritme vraagt tijd  $O(n)$ . Het product uit tijd en het aantal processoren is dus voor het parallelle algoritme  $O(n \log n)$  en voor het seriële  $O(n)$ . Dat betekent dat het algoritme niet kost optimaal is.

## Verbetering van de kost performantie

Wij zullen hier een trucje zien dat je soms kan gebruiken om de kost van een algoritme te verbeteren zonder de – asymptotische – tijd slechter te maken. Typisch voor gevallen waar dit trucje lukt is dat er maar één stap is waar alle processoren nodig zijn (een constant aantal stappen zou ook OK zijn). Eerst passen wij een licht gewijzigde versie van `copywaarden()` toe. Daarbij worden niet  $n$  kopieën gemaakt maar alleen  $n/\log n$  kopieën met  $O(n/\log n)$  processoren. De tijd die hiervoor nodig is, is  $O(\log(n/\log n)) = O(\log n - \log \log n) = O(\log n)$

Achteraf hebben wij ten minste  $n/\log n$  kopieën van  $a, n, \dots$ . Dan kopieert elke van de ongeveer  $n/(2 * \log n)$  processoren de waarden nog  $2(\log n - 1)$  keer. Dat vraagt nog eens tijd  $O(\log n)$  en achteraf hebben wij ten minste  $n$  kopieën, waarbij wij  $O(n/\log n)$  processoren gebruikt hebben en de tijd was  $O(\log n)$ . De kost is dus  $O((n/\log n) * \log n) = O(n)$  – wat dus identiek is aan het bestmogelijke seriële algoritme. Het algoritme is dus kost optimaal.

**Oefening 126** *Geef de pseudocode voor de processor met nummer `nummer` en het kost optimale algoritme.*

## Parallel lookup

Het volgende voorbeeld bepaalt in constante tijd of een getal  $a$  in een array aanwezig is:

**Algoritme 21** *Parallel Lookup in CRCW-PRAM*

```
is_present(a, lijst[])
// geeft TRUE terug als a in lijst[] -- anders FALSE

contained=FALSE

for i=1 to |lijst| pardo
    { if (lijst[i]=a) contained=TRUE; }

return contained;
```

Dit algoritme werkt alleen maar als je een CRCW-PRAM als model hebt. Alle  $|lijst|$  processoren moeten  $a$  tegelijk lezen en als er meerdere kopieën van  $a$  in `lijst` zitten, moeten ook meerdere processoren tegelijk schrijven. Deze pseudocode zegt duidelijk wat de  $|lijst|$  processoren in het parallel moeten doen – maar welke processor doet bv. `contained=FALSE`? Dit is

een deel dat van een arbitraire processor uitgevoerd kan worden – dus bv. processor  $p_1$ . Achteraf moeten alle processoren synchroon met hun deel in de lus beginnen en dan moet – nadat alle processoren klaar zijn – een arbitraire processor (dus bv. opnieuw  $p_1$ ) het commando **return contained** uitvoeren. Ook toekomstige pseudocode moet zo geïnterpreteerd worden.

**Algoritme 22** *Parallel Lookup in EREW-PRAM*

```
is_present(a, lijst)
{
  // geeft TRUE terug als a in lijst[] -- anders FALSE
  // n is de lengte van de lijst

  copywaarden(a,n,waarde[],aantal[])

  // ook n is nu gekend in aantal[]
  for i=1 to n pardo
    { contained[i]=FALSE;
      if (lijst[i]=waarde[i]) contained[i]=TRUE; }

  for j=1 to n/2 pardo
    { while (aantal[j]>1) // het aantal samen te vatten waarden
      aantal[j]=aantal[j]/2;
      if (j<= aantal[j])
      {
        contained[j]= contained[j] || contained[j+aantal[j]];
      }
    }
  // na deze lus staat TRUE in contained[1] als er ten minste 1
  // keer ergens TRUE stond

  return contained[1];
}
```

Je ziet dus dat het CR-model duidelijke voordelen heeft in vergelijking met een ER-model. Maar of dat in hardware ook zo geïmplementeerd kan worden is iets anders...

**Oefening 127** *Beschrijf een algoritme dat test of een getal  $p$  een priemgetal is. Het resultaat moet  $p$  zijn als  $p$  een priemgetal is en een niet triviale deler als  $p$  geen priemgetal is.*

- *Beschrijf een CRCW-PRAM algoritme dat in tijd  $O(1)$  werkt.*

- Beschrijf een EREW-PRAM algoritme dat in tijd  $O(\log p)$  werkt.

Wat is de inputlengte en hoeveel processoren gebruik je (als functie van de inputlengte)?

Hoeveel keer sneller is dit algoritme dan het priemzeef van Erathostenes op een seriële machine?

**Oefening 128** Beschrijf een algoritme voor een CRCW-machine dat in constante tijd het minimum van een verzameling van  $n$  getallen kan berekenen. Hoeveel processoren gebruik je?

**Tip:** Gebruik dat een getal het minimum is als alle mogelijke vergelijkingen met de andere getallen opleveren dat het kleiner is.

## 6.4 Semigroep algoritmen

Wij zullen nu een algoritme voor een EREW-machine beschrijven dat de som van  $n$  getallen in tijd  $O(\log n)$  kan berekenen.

Wij bespreken dat relatief abstract omdat wij zo in principe voor vele *verschillende* gevallen tegelijk een algoritme geven – bv. voor de operatoren “\*”, “minimum”, “maximum”, etc. Al deze operatoren hebben de eigenschap dat ze samen met de verzameling waarop ze opereren een semigroep vormen.

Een semigroep is een verzameling samen met een associatieve binaire bewerking – dus een bewerking met de eigenschap dat (als wij de bewerking als  $\oplus$  schrijven)  $(a \oplus b) \oplus c = a \oplus (b \oplus c)$ . Als bv. de verzameling de verzameling van gehele getallen is en de bewerking “+” dan voldoet dat zeker aan deze eigenschap. In feite heeft “+” zelfs nog de mooie eigenschappen dat de bewerking commutatief is en dat er een neutraal element en een invers element bestaan – maar dat hebben wij hier niet nodig. De voorbeelden zullen alleen maar voor verzamelingen van getallen zijn, maar als de operator  $\oplus$  bv. met matrices werkt, werken de algoritmen *in principe* op dezelfde manier.

De pseudocode voor deze algoritmen kan je als volgt schrijven:

**Algoritme 23** *Semigroep som in een EREW-PRAM*

//Invoer: array  $a[]$  met  $n = 2^k$  getallen beginnend met index 1.

//Resultaat:  $a[1] \oplus \dots \oplus a[n]$

```
semigroep_plus(a[],n)
```

```
{ copywaarden(n,n/2,n_array[],aantal[]);
```

```
  // nu kan elke processor die we gebruiken een kopie van n lezen
```

```

for j=1 to log n do
{
    for i=1 to n/2j pardo
        { buffer[i]=a[2*i-1]⊕a[2*i];
          a[i]=buffer[i]; }
    }

return a[1];
}

```

**Oefening 129** *Wijzig Algoritme 23 zo dat het ook voor getallen  $n$  werkt die geen macht van 2 zijn.*

Dit algoritme vraagt  $O(n)$  processoren en tijd  $O(\log n)$ . Omdat een serieel algoritme tijd  $O(n)$  vraagt is dit parallel algoritme dus niet kost optimaal – maar het is niet moeilijk om het zo te wijzigen dat het kost-optimaal is:

**Oefening 130** *Wijzig Algoritme 23 zo dat het de semigroep som in tijd  $O(\log n)$  op een kost optimale manier berekent.*

Als wij nu de werk complexiteit van Algoritme 23 berekenen dan zien wij dat elk deel van het algoritme inderdaad werk optimaal is: eerst wordt `copywaarden()` gebruikt en als wij daarvoor de kost optimale versie gebruiken dan werkt dat zeker met werk complexiteit  $O(n)$ .

Maar de for lus in `semigroep_plus()` vraagt werk

$$W(n) = c_1 + \sum_{j=1}^{\log n} (c_2 * n/2^j) = O(n)$$

met constanten  $c_1$  en  $c_2$ .

Dus vraagt ook dit deel van het Algoritme 23 werk  $O(n)$  en is dus werk optimaal.

**Maar:** Of deze analyse echt klopt, hangt er wel vanaf hoe de code van elke processor er precies uitziet. Als de processoren in de lus altijd toetsen of zij nog aan  $n/2^j$  voldoen, heeft elke van de processoren in alle  $\log n$  keren dat de buitenste lus wordt doorlopen werk. Dan is het werk  $O(n \log n)$ . Maar inderdaad kan een processor zodra hij niet meer aan deze voorwaarde voldoet helemaal stoppen – en dan is het werk in dit deel  $O(n)$ .

Als je het algoritme als een geheel ziet, moet je de  $O(n)$  processoren die de beschreven versie gebruikt al vanaf het begin hebben – dus ook al voor

copywaarden. Je moet de code voor die processoren dus zo schrijven dat zij het copywaarden gedeelte gewoon niet afwerken (een `if` gebaseerd op hun nummer dat dat deel overslaat) – en er dus ook geen werk doen.

**Oefening 131** (*Heel gemakkelijk – gewoon om te zien of de tekst helemaal verstaan is*)

*De volgende bewering heeft twee richtingen. Voor elke richting geef een bewijs of een tegenvoorbeeld:*

*Een parallel algoritme is werk optimaal als en slechts als het kost optimaal is.*

**Oefening 132** *Geef een  $O(\log n)$  tijdsbegrensd parallel EREW algoritme voor het berekenen van  $a[0] - a[1] + a[2] - a[3] + \dots + a[n-2] - a[n-1]$  voor een gegeven invoerarray  $a[]$  met  $n = 2^k \geq 2$  getallen beginnend met index 0.*

**Oefening 133** *Geef operatoren (en misschien een preprocessing dat constante tijd op een EREW-machine vraagt) zodat je het parallel lookup probleem voor een vast getal  $a$  (dat je dus in de code van elke processor kan schrijven) ook als een semigroep probleem kan beschouwen.*

## 6.5 Eén parallel grafenalgoritme

In dit deel zullen wij een parallel CRCW algoritme zien dat de samenhangscomponenten in een graaf bepaalt. Dat is aan de ene kant om te zien hoe parallel algoritmen op iets ingewikkeldere structuren dan alleen maar arrays van getallen werken en aan de andere kant om ook een algoritme te zien dat een beetje (maar niet veel) ingewikkelder is.

Seriële algoritmen voor het berekenen van samenhangscomponenten werken meestal met Depth-First search of Breadth-First search. Daarbij wordt aan elke top een getal toegekend zodat toppen die hetzelfde getal toegekend krijgen – het identificatienummer van de component – in dezelfde component zitten.

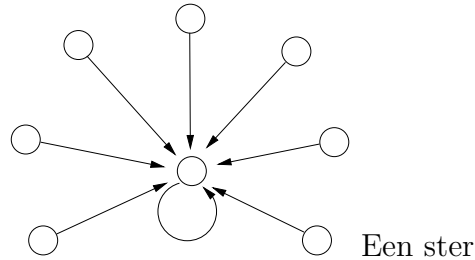
Hier willen wij ook dergelijke identificatienummers toekennen – maar dat in tijd  $O(\log |V|)$  als de graaf  $G = (V, E)$  is.

Wij gebruiken de union-find datastructuur die wij al in Datastructuren en Algoritmen II als voorbeeld hebben gezien om equivalentieklassen (dus ook samenhangscomponenten) efficiënt te berekenen en te updaten. Maar wij zullen de structuur hier opnieuw beschrijven.

**Definitie 8** *Een gerichte graaf  $D = (V, E)$  (met loops) heet een gerichte wortelboom als elke top één uitgaande boog heeft en er een  $v \in V$  (deze top*



noemen wij dan ook de wortel van de boom) bestaat waar de uitgaande boog een lus is en waar vanuit elke top een gericht pad naar  $v$  bestaat. Een gerichte wortelboom heet ster als de afstand van elke top tot de wortel ten hoogste 1 is.



**Opmerking 11** Stel dat in een CRCW PRAM een array `opv[]` van lengte  $n$  een gerichte graaf voorstelt waarin de enige cykels loops zijn en waarin uit elke top precies één boog vertrekt. Voor  $1 \leq i \leq n$  is `opv[i]` de eindtop van de unieke boog die in  $i$  vertrekt.

Dan kan een CRCW PRAM met  $n$  processoren in constante tijd beslissen welke top in een component ligt die een ster is – of precies: in constante tijd een array `ster[]` invullen zodat `ster[v]=TRUE` als  $v$  in een ster ligt en anders FALSE.

**Bewijs:** Het algoritme dat dit doet, werkt als volgt:

**Algoritme 24** *Bepalen welke toppen in sterren liggen*

```
// Het programma voor processor p
// In het begin: opv[] bevat de opvolgers
// Op het einde: ster[v] is juist ingevuld voor elke top v
```

```
void beleg_ster_array(ster[])
{
  ster[p]=TRUE;

  if (opv[p] != opv[opv[p]])
  { ster[p]=FALSE;
    ster[opv[opv[p]]]=FALSE;
  } // stap 1

  if (ster[opv[p]]==FALSE) ster[p]=FALSE; // stap 2
}
```

Dat dit algoritme in constante tijd werkt, is duidelijk. Het is ook duidelijk dat het hier kan gebeuren dat meerdere processoren dezelfde variabele willen lezen of daarin willen schrijven.

Dat dit algoritme juist werkt, kan als volgt gezien worden:

Wij noemen een top die zijn eigen opvolger is een looptop. Een top in een verzameling van wortelbomen maakt deel uit van een componente die een ster is als en slechts als hij aan één van de volgende voorwaarden voldoet:

- a.) hij is een looptop en er is geen gericht pad van lengte 2 naar deze top
- b.) hij is geen looptop, maar zijn opvolger voldoet aan a.)

Dat is misschien duidelijk, maar het is ook een goede oefening dat expliciet aan de hand van de definitie te bewijzen!

Na stap 1 staat voor elke looptop en elke top op een afstand van 2 of meer van een looptop de juiste waarde in `ster[]`. Een looptop die niet in een ster zit wordt door de processor  $p_v$  van een top  $v$  op afstand 2 als **FALSE** gemarkeerd en toppen  $w$  op afstand 2 of meer van een looptop worden door hun eigen processor  $p_w$  als **FALSE** gemarkeerd.

De enige toppen waarvoor de waarde op dit moment dus nog fout kan zijn, zijn toppen op afstand 1 van een looptop. Maar die maken natuurlijk precies dan deel uit van een ster als de looptop waar hun boog naartoe leidt deel uitmaakt van een ster – deze toppen worden dus in stap 2 juist gemarkeerd.



**Oefening 134** *Wordt de array `ster[]` ook door het volgende algoritme juist ingevuld?*

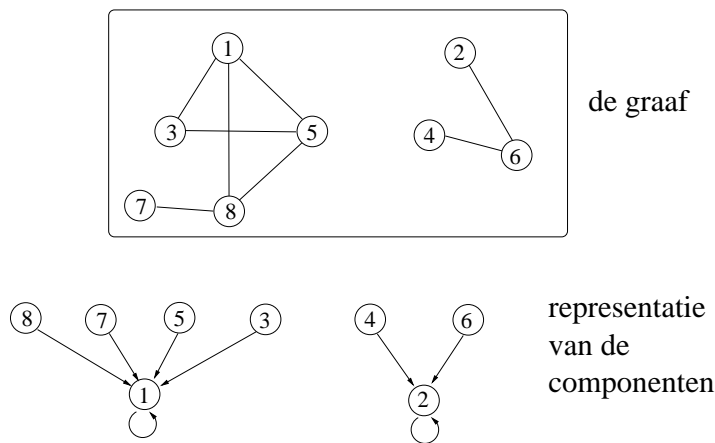
**Algoritme 25** `void beleg_ster_array(ster[])`

```
{
  ster[p]=TRUE;

  if (opv[p] != opv[opv[p]])
  { ster[p]=FALSE;
    ster[opv[p]]=FALSE;
    ster[opv[opv[p]]]=FALSE;
  }
}
```

}

De componenten worden voorgesteld door sterren – elke top wijst naar een top in de component – misschien zichzelf, misschien ook een andere top. Maar toppen in dezelfde component wijzen naar dezelfde top. Het nummer van de component waarin een top  $v$  zit, is dan gewoon de top waarnaar  $v$  wijst.



Figuur 20: Een graaf en hoe de samenhangscomponenten voorgesteld kunnen worden. Belangrijk is dat alle toppen in één component naar dezelfde top wijzen – het nummer dat de component voorstelt. Het is niet belangrijk welke van de toppen in de component dat is.

Het idee is nu dat in het begin alle toppen naar zichzelf wijzen. Deze singletons zijn de kleinstmogelijke wortelbomen. Die worden dan uitgebreid tot grotere wortelbomen en tegelijk wordt path compression (vergelijk data-structuren voor verzamelingen in DA2) toegepast. Ten slotte correspondeert elke component tot maar 1 wortelboom – en die is een ster of wordt tot een ster gecomprimeerd.

Het samenvoegen van de wortelbomen gebeurt in parallel. Voor elke boog  $\{v, w\}$  gebruiken wij een processor  $P_{\{v,w\}}$  die probeert de wortelboom waarin  $v$  zit en die waarin  $w$  zit samen te voegen (dit *samenvoegen* wordt ook *hooking* genoemd. Als de boom waarin  $v$  zit een ster is en die wordt aan die waarin  $w$  zit gehangen dan is de precieze *hooking* operatie  $\text{opv}[\text{opv}[v]] = \text{opv}[w]$ .

**Maar:** wij moeten opletten dat wij op deze manier geen gerichte cykels opbouwen door bv. tegelijk wortelboom  $B_1$  aan wortelboom  $B_2$  te voegen en omgekeerd! Bovendien moet de wortelboom met  $v$  een ster zijn omdat die anders door deze operatie uit elkaar kan vallen als  $\text{opv}[v]$  niet de wortel is.

Wij geven nu eerst de pseudocode en bespreken die dan in stappen. Om niet altijd met buffers te moeten werken die garanderen dat er in één lijn niet dezelfde variabele gelezen en geschreven wordt, spreken wij af, dat leesoperaties in één lijn altijd de waarde lezen van voor de lijn. Dat kan je dan ook met buffers schrijven, maar daardoor wordt de code alleen maar minder leesbaar.

**Algoritme 26** *Parallel samenhangscomponenten berekenen*

```
// Invoer: een graaf  $G=(V,E)$  als lijst van toppen en bogen
// Op het einde: voor elke top  $v$  bevat  $opv[v]$  een nummer dat
// zijn component kenmerkt.
```

```
// voor elke top  $v$  hebben wij een processor  $P_v$  en voor elke
// boog  $\{v,w\}$  een processor  $P_{\{v,w\}}$ 
// wij schrijven hook  $v \rightarrow w$  voor  $opv[opv[v]]=opv[w]$ 
```

```
mark_components(G) {
  for all vertices  $v$  pardo  $opv[v]=v$ ;
  // hier worden singletons gebouwd

  for all edges  $\{v,w\}$  pardo
  { if ( $v > w$ ) hook  $v \rightarrow w$ ; else hook  $w \rightarrow v$ ; //(a)
    if ( $v$  is singleton) hook  $v \rightarrow w$ ; //(b)
    if ( $w$  is singleton) hook  $w \rightarrow v$ ; //(b)
  }

  still_change=TRUE;
  while (still_change)
  {
    beleg_ster_array(ster[])
    for all edges  $\{v,w\}$  pardo
    { if ( $ster[v]$  AND ( $opv[v] > opv[w]$ )) hook  $v \rightarrow w$ ; //(c)
      if ( $ster[w]$  AND ( $opv[w] > opv[v]$ )) hook  $w \rightarrow v$ ; } //(c)
    beleg_ster_array(ster[]); // updaten
    for all edges  $\{v,w\}$  pardo
    { if ( $opv[v] \neq opv[w]$ )
      { if ( $ster[v]$ ) hook  $v \rightarrow w$ ; //(d)
        if ( $ster[w]$ ) hook  $w \rightarrow v$ ; } //(d)
      }
    // kijken of er veranderingen zijn en path compression:
    still_change=FALSE;
    for all vertices  $v$  pardo
```

```

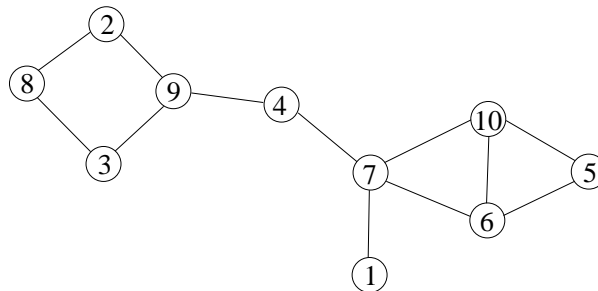
    { if (opv[v]≠opv[opv[v]])
      { opv[v]=opv[opv[v]]; still_change=TRUE; } //(e)
    }
  }
}

```

**Oefening 135** *In Algoritme 26 is niet precies beschreven hoe (v is singleton) in constante tijd beslist kan worden. Werk dit deel van de code expliciet uit. Herschrijf de hele parallelle lus die dit deel bevat.*

Nu zullen wij eerst naar een voorbeeld kijken hoe dit algoritme op een graaf werkt. In de tekeningen zijn de gewijzigde bogen altijd doorgaande lijnen terwijl oude bogen in de datastructuur met punten zijn voorgesteld. Natuurlijk is de voorgestelde ontwikkeling niet uniek. Als meerdere processoren in dezelfde variabele willen schrijven hangt het er vanaf welke processor erin slaagt te schrijven. Daarom staat in dergelijke gevallen aan de nieuwe bogen welke processor deze boog ingevuld heeft.

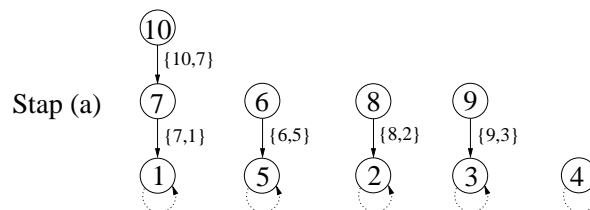
De graaf is:



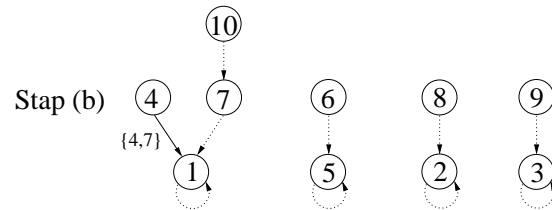
Eerst bouw je singletons:



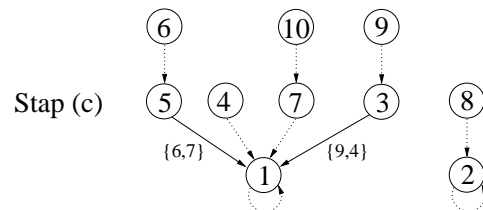
Dan bouw je de eerste wortelbomen in stap (a)



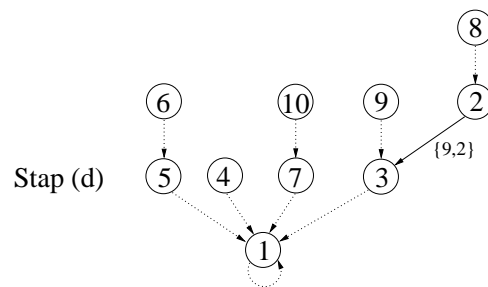
en hook je in stap (b) de overgebleven singletons aan bestaande wortelbomen:



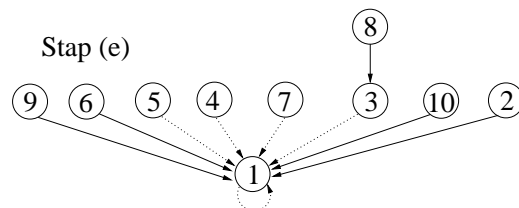
Nu begint de lus die sterren aan wortelbomen hookt en dan path compression toepast. Eerst worden in stap (c) sterren aan bestaande wortelbomen gehooked



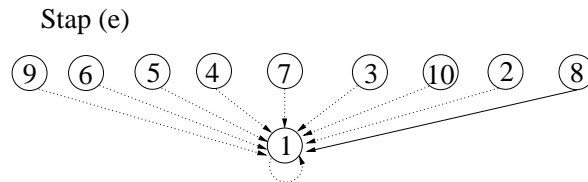
en dan in stap (d) nog eens overgebleven sterren:



Op het einde van de lus wordt path compression toegepast:



Er zijn nog wijzigingen, dus wordt de lus nog eens doorlopen. Er zijn geen sterren meer die nog gehooked kunnen worden, maar er is wel nog path compression mogelijk:



In de volgende doorloop van de lus zijn er dan helemaal geen wijzigingen meer en het programma stopt.

Om te bewijzen dat het algoritme juist werkt, moeten wij dus bewijzen dat er op het einde alleen maar sterren zijn en dat twee toppen dan en slechts dan op het einde tot dezelfde ster behoren als ze tot dezelfde component van de graaf behoren.

Dat bewijzen wij door verschillende kleine tussenstappen

**Opmerking 12** *Op elk moment vormt de datastructuur in het algoritme een woud van gerichte wortelbomen.*

**Bewijs:** In het begin zijn het allemaal singletons – daarvoor klopt de opmerking dus.

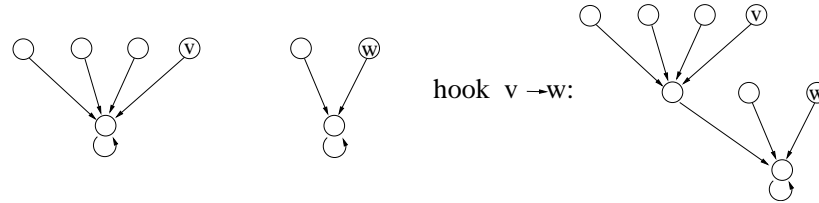
In stap (a) kunnen ook geen cyclen toegevoegd worden omdat alleen grotere toppen naar kleinere wijzen.

In stap (b) wordt alleen dan een boog vanuit  $v$  toegevoegd als alle burens van  $v$  groter zijn dan  $v$  (anders zou  $v$  nu een andere opvolger hebben) **en** al naar een andere top wijzen (de enige reden dat de poging van de corresponderende boog  $\{w, v\}$  niet slaagde, kan zijn dat de processor van een andere boog  $\{w, v'\}$  succesvol was). De burens zijn dus allemaal geen singletons meer – er wordt dus geen boog naar  $v$  toegevoegd en de datastructuur blijft dus een woud.

Na stap (b) zijn ook alleen nog die toppen singletons die ook in de graaf geïsoleerd zijn – alle andere toppen zijn deel van een wortelboom met ten minste 2 toppen.

De enige mogelijkheid dat in stap (c) een cykel geïntroduceerd wordt, is dat de bogen waarvan de processoren nieuwe hook-operaties uitvoeren componenten in een cykel aan elkaar voegen die allemaal sterren zijn (elke component moet een uitgaande boog hebben). Maar omdat het sterren zijn, worden in deze cykel altijd de wortels vergeleken als het om  $\text{opv}[v]$  en  $\text{opv}[w]$  gaat – en dus moet er één ster zijn met top  $v$  die gehooked wordt en waarbij  $\text{opv}[v] < \text{opv}[w]$  – een tegenstrijdigheid. Dus worden in deze stap ook geen cyclen gevormd.

In stap (c) worden geen nieuwe sterren gevormd omdat als één ster op een wortelboom gehooked wordt het resultaat in het begin een diepte van ten minste 2 heeft – zelfs dan als het allebei sterren zijn (zie Afbeelding 21).



Figuur 21: Een ster met top  $v$  wordt op een andere ster met top  $w$  gehooked.

Stap (d) heeft dus alleen sterren die ook voor stap (c) al sterren waren. Analooq met stap (b) worden ze nu alleen maar gehooked op componenten die geen sterren zijn waardoor opnieuw geen cykel kan ontstaan.

De enige stap waar `opv[]` dan nog gewijzigd wordt is stap (e). Maar het is duidelijk dat daar de paden alleen maar korter gemaakt worden en geen cyclen kunnen ontstaan.

■

**Opmerking 13** *Als een wortelboom voor de compressiestap (e) een diepte van  $d > 1$  heeft dan heeft hij achteraf een diepte van ten hoogste  $\lfloor (2d)/3 \rfloor$*

**Bewijs:** Met inductie kan je gemakkelijk bewijzen dat een top op even afstand  $d$  van de wortel voor de compressie achteraf een afstand van  $d' = d/2$  heeft. Als de afstand voor de compressie oneven is, is hij achteraf  $d' = (d + 1)/2$ . De maximale waarde van  $d'/d$  is dan voor  $d = 3$  en dan is  $d' = 4/2 = (2/3)d$ .

■

**Opmerking 14** *Als het algoritme stopt, zijn alle componenten sterren en twee toppen maken deel uit van dezelfde ster als en slechts als ze tot dezelfde component behoren.*

**Bewijs:** Zolang er nog componenten zijn die geen sterren zijn, wordt doorgegaan met de compressie en dus zijn er ook wijzigingen en het algoritme gaat door. Dat nooit toppen uit verschillende componenten in dezelfde wortelboom terechtkomen is ook duidelijk. Stel nu dat het algoritme



klaar is en dat er nog twee toppen zijn die tot dezelfde component behoren maar in verschillende wortelbomen zitten – dus in verschillende sterren. Deze toppen kunnen zo gekozen worden dat er een boog  $\{v, w\}$  is die beide toppen bevat (waarom?). Wij mogen stellen dat  $opv[v] > opv[w]$  – maar dan worden de wortelbomen van  $v$  en  $w$  ofwel aan elkaar gehooked ofwel aan andere wortelbomen – maar het algoritme is zeker nog niet afgelopen – een tegenstrijdigheid.

■

Wij weten dus dat het algoritme juist werkt en moeten alleen nog de complexiteit bepalen:

**Lemma 15** *De tijdscomplexiteit van Algoritme 26 opgestart op een graaf  $G = (V, E)$  is  $O(\log |V|)$  met  $|V| + |E|$  processoren.*

**Bewijs:** Tot de **while** (**still\_change**) lus is zeker maar constante tijd nodig en het is ook duidelijk dat één doorloop van de lus constante tijd vraagt. Wat nog bewezen moet worden is dat de lus maar  $O(\log |V|)$  keren wordt doorlopen.

Daarvoor definiëren wij een wortelboom in de datastructuur als *actief* (in een doorloop van de lus) als hij in de lus op een andere wortelboom gehooked wordt, een andere wortelboom op deze gehooked wordt of als  $w$  gecomprimeerd wordt. Anders heet de boom *passief*. Een boom is passief als en slechts als het een ster is en er is geen boog van een top in deze ster naar een top in een andere wortelboom: als een wortelboom geen ster is wordt hij zeker gecomprimeerd en als hij een ster is met burens in andere wortelbomen dan wordt hij zeker in stap (c) of (d) gehooked. Maar dat betekent ook dat als een wortelboom één keer passief is hij dat zeker ook blijft!

Nu kijken wij naar de volgende waarde van de datastructuur  $D$ :

$$W(D) = \sum_{\{T \in D \mid T \text{ is actief}\}} d(T)$$

waarbij  $d(T)$  de diepte van de wortelboom  $T$  is. Het is belangrijk te zien dat omdat een wortel nooit aan een blad geplakt wordt maar de pointer naar een top op afstand ten minste 1 van een blad wijst, de diepte van een verzameling wortelbomen die op elkaar gehooked wordt ten hoogste de som van de enkele diepten voor het hooken kan zijn. Tijdens de hooking operaties kan  $W(D)$  dus niet groeien. Inderdaad kan  $W(D)$

dus alleen maar kleiner worden en als path compression wordt toegepast of componenten pasief worden zal dat ook zeker gebeuren.

Natuurlijk geldt altijd  $W(D) \leq |V|$  en als  $W(D_t)$  de waarde op het einde van de  $t$ -de doorloop van de lus is dan geldt met Opmerking 13  $W(D_{t+1}) \leq (2/3)W(D_t)$ . Stel dat  $t$  het nummer van de op één na laatste doorloop is. Dan geldt  $W(D_t) \geq 1$  omdat het algoritme anders zou stoppen (geen wijzigingen). Als  $D_0$  de datastructuur voor de eerste doorloop van de lus is dan geldt met  $W(D_0) \leq |V|$ :

$$W(D_t) \leq \left(\frac{2}{3}\right)^t W(D_0) \Rightarrow 1 \leq \left(\frac{2}{3}\right)^t |V| \Rightarrow t \leq \frac{1}{\log(3/2)} \log |V|$$

Dus wordt de lus inderdaad maar  $O(\log |V|)$  keren doorlopen.

■

**Oefening 136** *Wijzig het algoritme zo, dat het aantal componenten van een graaf bepaald en uitgevoerd wordt.*

**Oefening 137** *Wijzig het algoritme zo, dat de nummers van de componenten niet arbitraire toppennummers zijn maar  $1, \dots, k$  als er  $k$  componenten zijn. (Deze oefening is niet zo gemakkelijk.)*

**Oefening 138** *Bewijs expliciet dat twee toppen die in Algoritme 26 in dezelfde wortelboom terechtkomen ook tot dezelfde component van de graaf behoren.*

**Oefening 139** *Gegeven een verzameling  $M$  van elementen (in DA2 hebben wij dat altijd een universum genoemd) en een verzameling  $R$  van relaties tussen deze elementen. Beschrijf een parallel CRCW algoritme voor een PRAM dat de door  $R$  gegenereerde equivalentieklassen in tijd  $O(\log |M|)$  berekent.*

**Oefening 140** *Bespreek de wijzigingen die in Algoritme 26 nodig zijn zodat het ook op een EREW PRAM kan werken. Wat is de complexiteit van jouw gewijzigd algoritme?*

## 6.6 Pipelining

Ten slotte zullen wij het nog kort over een netwerkmodel hebben – alleen om een dergelijk model ook eens gezien te hebben. In dit model werken wij met gewone processoren die allemaal hun eigen geheugen hebben. Als in de pseudocode dus een variabele gebruikt wordt dan is die lokaal – dus verschillend

voor verschillende processoren. Maar ze kunnen ook data doorsturen naar andere processoren waarmee ze verbonden zijn. Wij gebruiken bevelen zoals `send(x,i)` of `receive(y,j)` om aan te duiden dat wij de inhoud van een variabele naar processor  $i$  sturen of een boodschap van processor  $j$  in variabele  $y$  willen opslaan. Als op processor  $j$  het bevel `send(x,i)` gebruikt wordt, moet er wel een verbinding zijn tussen de processoren  $i$  en  $j$  – anders moet ervoor gezorgd worden dat de data met hulp van andere processoren doorgestuurd wordt – wat natuurlijk duurder is. In ons voorbeeld is het voldoende als voor  $n$  processoren  $1, \dots, n$  processor  $i$  voor  $1 \leq i \leq n - 1$  een lees/schrijf verbinding heeft met processor  $i + 1$ . Het netwerk is hier dus heel eenvoudig – het is een pad.

Zoals de naam pipelining al doet vermoeden geven deze algoritmen data door langs een ketting van processoren. En ook al is het algoritme ontworpen voor een netwerk van processoren kan het op een PRAM natuurlijk gemakkelijk gesimuleerd worden. In feite kunnen deze pipelining algoritmen vooral op een multicore processor met linux gebruik makend van pipes bijzonder gemakkelijk geïmplementeerd worden.

Een kenmerk dat aantoont dat het principe van pipelining toepasbaar is, is een soort genestelde structuur van de problemen. Het principe kan het best op voorbeelden gedemonstreerd worden:

Het eerste voorbeeld dat wij zullen zien, is het sorteren van getallen. Invoer zijn  $n$  getallen (die van processor 1 worden gelezen) en uitvoer zijn deze  $n$  getallen in dalende volgorde die van processor 1 worden geschreven. Als wij het grootste van de  $n$  getallen verwijderen dan is het tweede grootste het grootste van de overblijvende getallen – en algemeen is het  $i$ -de grootste het grootste van de verzameling van getallen waaruit de grootste  $i - 1$  getallen al verwijderd zijn.

Dit kunnen wij als volgt implementeren: processor 1 leest alle getallen, houdt het grootste getal bij en stuurt de rest naar processor 2. Algemeen houdt processor  $i$  ( $2 \leq i \leq n$ ) het grootste van de getallen die het van processor  $i - 1$  krijgt bij en stuurt de rest naar processor  $i + 1$ . Als wij `receive(.,0)` interpreteren als *lees het volgende getal uit de invoer* en `send(.,0)` als *schrijf het volgende getal naar de uitvoer* en het einde van de invoer door een EOF gekenmerkt is, kunnen wij dit algoritme opschrijven als

#### **Algoritme 27** *Sorteren met een pipeline*

```
Het programma voor processor p
//Invoer: een reeks van n getallen
//Uitvoer: dezelfde reeks in dalende volgorde

{
```

```

receive(best,p-1);

if (best!=EOF)
{
    receive(temp,p-1);
    while (temp != EOF)
        { if (temp>best)
            { send(best,p+1); best=temp; }
          else send(temp,p+1);
          receive(temp,p-1);
        }
    // nu zijn alle getallen gelezen
    send(EOF,p+1);

    send(best,p-1);
    // het grootste van de ontvangen getallen wordt eerst gestuurd
    receive(temp2,p+1);
    while (temp2 != EOF)
        {
            send(temp2,p-1);
            receive(temp2,p+1);
        }
    }
send(EOF,p-1);
}

```

Op het einde worden dus de getallen aan de vorige processor teruggestuurd – en omdat altijd het grootste getal eerst wordt gestuurd, kan je gemakkelijk met inductie bewijzen dat de getallen in dalende volgorde uitgevoerd worden.

Wij zullen nu nog een tweede voorbeeld zien waarvan jullie het seriële algoritme al kennen: het priemzeef van Eratosthenes.

Invoer zijn hier alle getallen  $2, 3, \dots, n$  en uitvoer zijn alle priemgetallen tussen 2 en  $n$ . De genestelde structuur kan hier zo beschreven worden dat een getal  $x$  een priemgetal is als en slechts als het door geen priemgetal kleiner dan  $x$  deelbaar is. De  $i$ -de processor in de rij krijgt allen maar getallen die niet door de eerste  $i - 1$  priemgetallen deelbaar zijn en schrijft alleen die, die niet door de eerste  $i$  priemgetallen deelbaar zijn – verwijdert dus die getallen die door *zijn* priemgetal deelbaar zijn. Daarbij is *zijn priemgetal* het kleinste getal dat deze processor krijgt – en omdat de getallen in volgorde worden doorgegeven is het dus het eerste getal. Als wij opnieuw processor nummer 0 met invoer en uitvoer identificeren dan is het programma voor processor  $p$ :

**Algoritme 28** *Priemzeef met een pipeline*

```
//Het programma voor processor p
//Invoer: de getallen 2,...,n
//Uitvoer: de priemgetallen 2 <= p <= n in stijgende volgorde

{
receive(getal,p-1);
temp=getal;

while (temp != EOF)
{ if (temp modulo getal !=0) send(temp, p+1);
  receive(temp,p-1);
}
// nu zijn alle getallen gelezen

send(EOF, p+1);

if (getal!=EOF)
{ send(getal,p-1);
  while (receive(temp,p+1)!=EOF) send(temp,p-1);
}
send(EOF,p-1);
}
```

Ook hier kan je door middel van inductie gemakkelijk bewijzen dat dit algoritme juist werkt.

Pipelining is in feite meer dan alleen maar toepasbaar op dit eenvoudige netwerkmodel. Het is een techniek die ook voor de ontwikkeling van algoritmen voor PRAMs toegepast kan worden als een taak mooi in deeltaken opgesplitst kan worden die na elkaar uitgevoerd moeten worden. In feite zijn op een PRAM veel dingen nog gemakkelijker omdat je natuurlijk een gemeenschappelijk geheugen kan gebruiken. Maar dit opsplitsen in deeltaken die **na** elkaar uitgevoerd moeten worden heeft natuurlijk al iets dat *inherent serieel* is en soms is de speedup niet zo groot als je zou hopen of kan op een PRAM ook door andere middelen verkregen worden.

**Oefening 141** *Gegeven een functie  $f : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  waarvan je mag veronderstellen dat ze al geïmplementeerd is en op een enkele processor in constante tijd berekend kan worden. Invoer is nu een reeks  $a_1, \dots, a_n, -, b_1, \dots, b_n$  en het doel is nu een functie*

$$\sum_{i=1}^n \sum_{j=1}^n f(a_i, b_j)$$

te berekenen. Dat moet dus de uitvoer zijn.

Geef de pseudocode voor een pipeline algoritme dat dit resultaat in lineaire tijd berekent. Invoer en uitvoer gebeuren zoals in de voorbeelden vanuit processor 1.

**Oefening 142** Beschrijf een algoritme voor een CRCW-PRAM met  $n^2$  processoren dat  $n$  verschillende getallen in een array  $a[]$  in tijd  $O(\log n)$  in dalende volgorde in de array  $s[]$  kan schrijven.

## Index

- $K(n)$ , 139
- $N(a)$ , 6
- $W()$ , 139
- adaptieve Huffman codering, 114
- afgewerkt
  - MPI\_Send(), 154
- afstand
  - Levenshtein, 95
- algoritme
  - Boyer-Moore-Horspoolshift, 85
  - genetisch, 26, 31
  - Knuth-Morris-Pratt, 75
  - parallel, 136
  - Rabin-Karp, 71
  - shift-AND, 89
  - shift-AND voor matches met fouten, 93
- Baeza-Yates, R., 89
- blokkeren, 154
- Bloom filter, 48
- Boyer, 85
- Boyer-Moore-Horspool, 85
- Breadth-First search, 168
- buren, 6
- Burrows-Wheeler transformatie, 127
- Burton Bloom, 48
- bzip2, 127
- communicator, 147
- computing
  - distributed, 138
  - parallel, 136
- Concurrent Read, 161
- Concurrent Write, 161
- CR, 161
- crossover, 26, 28, 29
- CW, 161
- datatypes
  - MPI, 150
- deadlock, 154
- deletion, 94
- Depth-First search, 168
- distributed computing, 138
- doelfunctie, 3
- dynamisch programmeren, 98
- editeerafstand, 95
- ER, 161
- Eratosthenes
  - priemzeef van, 180
- EW, 161
- Exclusive Read, 161
- Exclusive Write, 161
- extendible hashing, 38
- false positives, 49
- fenotype, 27
- filter
  - Bloom, 48
- fitness, 26
- GA, 27, 31
- geleid lokaal zoeken, 12
- genetisch algoritme, 31
- genetische algoritmen, 26
- genotype, 27
- gesimuleerd temperen, 22
- Gonnet, G., 89
- grafsteen, 59
- guided local search, 12
- hooking, 171
- Horspool, 85
- Huffman codering, 110
  - adaptief, 114
- karakteristieke vector, 90

- Karp, 71
- Knuth, D., 75
- kost, 139
- kost optimaal, 139
  
- laadfactor, 43
- Levenshtein, 95
- Local search, 5
- local search, 6
- lokaal zoeken, 6
  
- Manber, U., 93
- Message Passing Interface, 145
- Metaheuristiek, 3
- Moore, 85
- Morris, J.H., 75
- move to front, 134
- MPI, 145
- MPI datatypes, 150
- MPI handle, 156
- MPI\_Comm\_rank(), 148
- MPI\_Comm\_size(), 148
- MPI\_Finalize(), 146
- MPI\_Init(), 146
- MPI\_Irecv(), 156
- MPI\_Isend(), 156
- MPI\_Recv(), 151
- MPI\_Send(), 150
- MPI\_SUCCESS, 147
- MPI\_Test(), 157
- mpicc, 146
- mpirun, 146
- mutatie, 28
- mutaties, 26
  
- niveau, 141
  
- objective function, 3
- overstroomemmers, 43
  
- parallel computing, 136
- parallele algoritmen, 136
  
- path compression, 171
- pipelining, 179
- point-to-point communicatie, 149
- PRAM, 160
- Pratt, V., 75
- priemzeef van Eratosthenes, 180
  
- Rabin, 71
- receive(y,j), 179
- root proces, 148
  
- SA, 24
- selectie, 31
- semigroep, 166
- send(x,i), 179
- shift-AND, 89
  - matches met fouten, 93
- simulated annealing, 22, 24
- singletons, 171
- ster, 169
  
- toevallig zoeken, 5
- transformatie
  - Burrows-Wheeler, 127
- transputers, 137
  
- variable neighbourhood descent, 18
- variable neighbourhood search, 18, 20
- vector
  - karakteristiek, 90
- verschuivingstabel, 76
- vnd, 18
- vns, 20
  
- werk complexiteit, 139
- werk optimaal, 139
- wortelboom, 168
  - gericht, 168
- Wu, S., 93