



SAPIENZA  
UNIVERSITÀ DI ROMA

## Calcolo dei numeri reali esatti in Haskell

Facoltà di Ingegneria dell'informazione, informatica e statistica  
Corso di Laurea in Informatica

Candidato

Andrea Princic

Matricola 1837592

Relatore

Prof. Ivano Salvo

Anno Accademico 2020/2021

Tesi discussa il 14 Dicembre 2021  
di fronte a una commissione esaminatrice composta da:  
Enrico Tronci (presidente)  
Salvatore Pontarelli  
Gabriele Tolomei  
Adolfo Piperno  
Giuseppe Pirrò  
Daniele Gorla  
Ivano Salvo

---

**Calcolo dei numeri reali esatti in Haskell**

Tesi di Laurea. Sapienza – Università di Roma

© 2021 Andrea Princic. Tutti i diritti riservati

Questa tesi è stata composta con L<sup>A</sup>T<sub>E</sub>X e la classe Sapthesis.

Versione: 9 dicembre 2021

Email dell'autore: princic.1837592@studenti.uniroma1.it

## Sommario

Il tema della tesi è il calcolo dei numeri reali esatti, ovvero la possibilità di eseguire operazioni usando numeri in notazione posizionale con un numero di cifre potenzialmente infinito e non periodico con precisione arbitraria. Questo è un problema importante in applicazioni pratiche che necessitano di fare calcoli molto precisi nei quali un piccolo errore di approssimazione potrebbe propagarsi e ingigantirsi nel susseguirsi dei calcoli, andando a impattare sulla precisione del risultato finale della computazione.

# Indice

<b>1</b>	<b>Computazione esatta sui numeri reali</b>	<b>2</b>
1.1	Il problema . . . . .	2
1.2	Numeri computabili . . . . .	3
1.3	Rappresentazioni per i reali . . . . .	4
1.3.1	Gli stream . . . . .	4
1.3.2	Stream di cifre . . . . .	5
1.3.3	Cifre binarie con segno . . . . .	6
1.3.4	Basi non naturali . . . . .	8
1.3.5	Rappresentazione in base $\phi$ . . . . .	9
<b>2</b>	<b>La Golden Notation</b>	<b>13</b>
2.1	Addizione semplificata . . . . .	14
2.2	Addizione completa . . . . .	14
2.3	Complemento semplificato . . . . .	15
2.4	Complemento e Sottrazione completi . . . . .	16
2.5	Prodotto semplificato . . . . .	16
2.6	Prodotto completo . . . . .	17
2.7	Divisione semplificata . . . . .	18
2.8	Divisione completa . . . . .	18
<b>3</b>	<b>Implementazione</b>	<b>19</b>
3.1	Il linguaggio Haskell . . . . .	19
3.2	Tipi per la rappresentazione . . . . .	19
3.3	Prodotto completo . . . . .	20
3.4	Costanti . . . . .	21
3.5	Conversione . . . . .	23
3.5.1	Conversione di un intero . . . . .	23
3.5.2	Conversione di un razionale . . . . .	24
3.6	Fibonacci . . . . .	25
3.6.1	Senza ricorsione . . . . .	25
3.6.2	Con ricorsione . . . . .	25
	<b>Bibliografia</b>	<b>27</b>

# Introduzione

La computazione esatta fa uso di rappresentazioni e algoritmi tali che si possa raggiungere una precisione arbitraria prima, durante e dopo l'esecuzione di una serie di operazioni sui numeri reali. Queste rappresentazioni possono essere molto diverse da quelle usate comunemente nei calcolatori e fanno uso di strutture dati virtualmente infinite, chiamate stream. L'utilizzo di queste rappresentazioni e rispettivi algoritmi potrebbe comportare una perdita dal punto di vista dell'efficienza nell'esecuzione delle operazioni, ma offre il vantaggio di non avere perdite di precisione durante tutto il processo di calcolo.

La tesi è suddivisa in 3 parti fondamentali:

- Computazione esatta sui numeri reali: in questo capitolo si introduce il problema del calcolo esatto di numeri reali e la principale rappresentazione usata nei moderni calcolatori per i numeri razionali, e il motivo per cui non è adatta a fare calcoli esatti. Si illustrano anche alcune possibili soluzioni per rappresentare e calcolare numeri reali a precisione arbitraria. Queste rappresentazioni introducono una ridondanza per rendere computabili le operazioni sui reali utilizzando una quantità finita di memoria.
- La Golden Notation: in questo capitolo si illustrano la rappresentazione e gli algoritmi scelti per questa tesi. La base per la Golden Notation è il numero reale  $\phi$ , scelto per la semplicità con cui si definiscono le operazioni aritmetiche. In questa base la ridondanza è data dall'uguaglianza  $1.00 = 0.11 = 0.\overline{10}$ .
- Implementazione: in questo capitolo si discute l'implementazione dei suddetti rappresentazione e algoritmi nel linguaggio di programmazione puramente funzionale Haskell. Oltre ad alcune operazioni della Golden Notation, vengono illustrate le funzioni per il calcolo dei numeri di Fibonacci e un algoritmo per convertire un numero razionale di grandezza illimitata in base 10 in un numero in Golden Notation.

## Capitolo 1

# Computazione esatta sui numeri reali

### 1.1 Il problema

Il problema della rappresentazione dei numeri reali all'interno dei calcolatori elettronici è un problema che riguarda principalmente la precisione. Alcuni insiemi di numeri ( $\mathbb{N}$ ,  $\mathbb{Z}$ ) sono facilmente rappresentabili nel sistema binario come sequenze di bit. Questi, nella rappresentazione classica con un numero fisso di bit, in base al numero di bit usati per rappresentarli hanno un limite superiore di grandezza (e uno inferiore per i numeri con segno), ma rimangono comunque numeri precisi, senza approssimazioni. Per questi due insiemi esistono anche rappresentazioni con un numero variabile di bit, che permettono di rappresentare naturali e interi di grandezza illimitata grazie ad un incremento nel numero dei bit utilizzati nella rappresentazione, quando necessario. Esempi di questa rappresentazione sono gli interi in Python e Haskell.

Non vale la stessa cosa per i numeri razionali ( $\mathbb{Q}$ ) né tanto meno per i numeri reali ( $\mathbb{R}$ )

I numeri razionali vengono solitamente rappresentati con una rappresentazione chiamata mantissa-esponente, dal significato simile alla notazione scientifica, che rappresenta un numero come una coppia composta da:

- mantissa: una sequenza di bit, con segno positivo o negativo, la cui lunghezza determina la precisione del numero
- esponente: un numero intero che rappresenta l'effettiva grandezza del numero rappresentato

Usando questa rappresentazione abbiamo che una coppia mantissa-esponente  $(m, e)$  nel sistema binario rappresenta il numero:

$$\frac{m}{2^{p-1}} \times 2^e$$

dove  $p$  è il numero di bit della mantissa. Questa rappresentazione può ovviamente essere usata con qualunque altra base diversa da 2.

In qualunque rappresentazione a base naturale (ad esempio decimale o binaria) i numeri razionali hanno la proprietà di poter essere rappresentati con una quantità finita di memoria. Più aumenta la precisione della rappresentazione più memoria sarà necessario usare ma, comunque, rimarrà sempre una quantità finita. Questo, però, non accade con i numeri reali non razionali. I numeri reali non sono rappresentabili con una quantità finita di memoria (o di cifre) in quanto la loro parte decimale è infinita (come anche quella dei razionali) ma ha la particolarità di non essere periodica. Questo implica che per rappresentarli servirebbe una quantità infinita di memoria, una cosa impossibile sugli attuali calcolatori elettronici.

Dunque i numeri rappresentabili in un moderno calcolatore sono soltanto un sottoinsieme finito dei numeri razionali, la cui grandezza dipende dalla quantità di memoria destinata a rappresentarli. Qualunque numero che non appartenga a quel sottoinsieme rappresentabile viene necessariamente arrotondato al più vicino numero rappresentabile, causando un'inevitabile perdita di precisione sul risultato finale.

Un esempio di questa imprecisione è il risultato della somma tra i numeri razionali 0.1 e 0.2 fatta con il tipo float (32 bit):

$$0.1f + 0.2f = 0.300000011920928955078125$$

e con il tipo double (64 bit):

$$0.1d + 0.2d = 0.3000000000000000444089209850062616169452667236328125$$

Oppure, provando ad usare il tipo float sul più grande valore rappresentabile da un intero a 32 bit: 2147483647 si ottiene invece il valore 2147483648 e facendovi alcune operazioni questi sono i risultati:

$$\begin{aligned} 2147483647.0f &= 2147483648.000000 \\ 2147483647.0f - 1 &= 2147483648.000000 \\ 2147483647.0f - 64 &= 2147483648.000000 \\ 2147483647.0f - 65 &= 2147483520.000000 \end{aligned}$$

Questo accade perché a quella precisione il range che separa due float adiacenti è 128, quindi per ottenere un cambiamento bisogna sottrarre (o sommare) almeno la metà di questo range, passando così al float precedente (o successivo) [2].

Nei prossimi capitoli vengono esposte alcune rappresentazioni che permettono di risolvere questo problema.

## 1.2 Numeri computabili

I numeri reali che si possono calcolare con le tecniche che seguono non sono tutti gli elementi di  $\mathbb{R}$ , anzi, sono soltanto un sottoinsieme di tutti i numeri reali. In particolare, i numeri reali computabili sono un sottoinsieme numerabile dell'insieme  $\mathbb{R}$ . Dal momento che  $\mathbb{R}$  ha una cardinalità maggiore di  $\mathbb{N}$ , questo significa che esistono

dei numeri reali che non si possono conoscere, né calcolandoli con un programma né con qualche formula matematica, a nessun livello di approssimazione.

Esistono tante definizioni dell'insieme dei numeri computabili, tutte equivalenti. Una di queste [3] utilizza le macchine di Turing per definire un numero computabile come un numero  $x$  per il quale esiste una macchina di Turing tale che, dato un numero  $n$  sul nastro iniziale, termina con l' $n$ -esima cifra di  $x$  sul nastro. Da questa definizione si può dedurre che, essendo le macchine di Turing numerabili, anche i numeri computabili sono al massimo numerabili.

Un'altra definizione è quella che definisce computabile un numero  $a$  se può essere approssimato da qualche funzione computabile  $f : \mathbb{N} \rightarrow \mathbb{Z}$  in questo modo: dato un numero naturale  $n$ , la funzione produce un intero  $f(n)$  tale che:

$$\frac{f(n) - 1}{n} \leq a \leq \frac{f(n) + 1}{n}$$

Una funzione computabile è una funzione che può essere computata da una macchina di Turing o da un modello di computazione equivalente.

## 1.3 Rappresentazioni per i reali

Come fare, dunque, per rappresentare un numero reale, o per lo meno un razionale a precisione arbitraria, in un calcolatore elettronico? Prima di tutto servirebbe una struttura dati virtualmente infinita, qualcosa che permetta di immagazzinare una quantità non finita di dati, ovviamente non tutti insieme (la memoria di un calcolatore è pur sempre finita). Una struttura che si presta a questo scopo è lo stream.

### 1.3.1 Gli stream

Gli stream non sono vere e proprie strutture dati. In generale, uno stream è una sequenza potenzialmente infinita di elementi, i quali però non vengono elaborati tutti insieme ma soltanto uno alla volta, quando è necessario. Gli stream quindi si prestano bene allo scopo di rappresentare sequenze potenzialmente infinite di cifre, dato che per calcolarne una parte non è necessario calcolarle per forza tutte.

Si potrebbe ad esempio ipotizzare che in uno stream infinito sia rappresentato un numero reale, con tutte le sue infinite cifre. Ora, anche se non se ne può ottenere il valore esatto, si può comunque ottenere una quantità arbitrariamente grande di cifre, tale da rendere il valore ottenuto preciso a tal punto da poterne trascurare l'inesattezza.

Gli algoritmi che utilizzano gli stream, non potendo utilizzare tutto il loro contenuto insieme, fanno tipicamente uso di alcune cifre di lookahead, ovvero utilizzano soltanto un numero finito di cifre dalla testa degli stream in input per generare una o più cifre dello stream in output, per poi continuare ricorsivamente a generare cifre utilizzando il resto dello stream. Questi algoritmi quindi funzionano in modo molto diverso dai classici algoritmi per le operazioni a cui siamo abituati, che iniziano facendo i conti dall'ultima cifra a destra e vanno verso sinistra. Gli algoritmi sugli stream non possono accedere all'ultima cifra (in quanto non esiste un'ultima cifra) quindi eseguono i calcoli a partire da sinistra e procedono verso destra.



Nei prossimi capitoli si useranno le seguenti sintassi per riferirsi ad uno stream:

$$[x_1, x_2, x_3, \dots] \equiv x_1 : x_2 : x_3 : \dots$$

La seguente sintassi per riferirsi alla concatenazione di cifre in testa ad uno stream, dove  $x$  è uno stream e  $x_i$  è una cifra:

$$x_1 : x_2 : x_3 : x$$

La seguente sintassi per riferirsi ad uno stream infinito della cifra  $x$

$$\overrightarrow{x}$$

Inoltre va fatta una distinzione tra il concetto di numero e di numerale: un numerale è un elemento sintattico mentre un numero è la sua parte semantica. Un numero è il valore semantico che si attribuisce a un numerale dopo la sua interpretazione.

Esempi di numerali sono gli stream e le coppie mantissa-esponente. Per differenziare questi due concetti si usano le doppie parentesi quadre  $\llbracket x \rrbracket$  per indicare il numero rappresentato dal numerale  $x$ .

Un'ultima cosa da notare è che, nella maggior parte dei casi, nei prossimi capitoli si useranno gli stream per rappresentare soltanto un intervallo chiuso di numeri reali e non l'intera linea, visto che negli stream manca l'informazione che servirebbe a distinguere la parte intera dalla parte decimale di un numero. In generale, uno stream

$$x = d_1 : d_2 : d_3 : \dots$$

in base  $b$  verrà interpretato come

$$\llbracket x \rrbracket = \sum_{i=1}^{\infty} d_i \times b^{-i} = \frac{d_1}{b} + \frac{d_2}{b^2} + \frac{d_3}{b^3} + \dots$$

### 1.3.2 Stream di cifre

Una prima e semplice rappresentazione dei numeri reali nell'intervallo  $[0, 1]$  tramite stream potrebbe essere quella di rappresentarli come semplici sequenze di cifre, decimali o binarie che sia, da sinistra a destra. Questa rappresentazione, purtroppo, non può funzionare.

Si prenda ad esempio il caso dei numeri  $0.333\dots$  e  $0.666\dots$  in base 10, rappresentati dai seguenti stream:

$$3 : 3 : 3 : \dots \qquad 6 : 6 : 6 : \dots$$

e si ipotizzi di volerne fare la somma. Tale operazione non sarebbe computabile su questi due numeri in questa rappresentazione perché per determinare anche soltanto la prima cifra del risultato servirebbe una quantità infinita di input letto dai due addendi.

Questo accade per via del fatto che uno stream che inizia con una lunga sequenza di 3 non per forza rappresenta il numero  $\frac{1}{3}$ , e allo stesso modo uno stream che inizia con una lunga sequenza di 6 non per forza rappresenta il numero  $\frac{2}{3}$ . Potrebbe essere

che ad un certo punto di quegli stream, magari dopo migliaia di cifre, si trovi una cifra diversa dal 3 o dal 6, ad esempio la cifra 0. In tal caso il risultato della somma non sarebbe 1 ma un numero che si trova nell'intervallo

$$[0.9, 1)$$

ma per stabilirlo servirebbe analizzare una quantità potenzialmente infinita di cifre da entrambi gli stream.

Lo stesso problema sorge anche nel caso di basi diverse da 10: ad esempio in base 2 si prendano in esame gli stream

$$1 : 1 : 1 : \dots \quad \text{e} \quad 0 : 0 : 0 : \dots$$

che potrebbero rappresentare i numeri 1 e 0, ma non è possibile saperlo senza esaminarne potenzialmente tutte le cifre.

Questa rappresentazione è quindi da escludere dalle possibili soluzioni.

### 1.3.3 Cifre binarie con segno

La seconda rappresentazione [4] è simile alla prima (ora si parla di base 2) ma aggiunge una nuova cifra alle due usate normalmente: la cifra  $-1$  alla quale successivamente, nella sintassi degli stream, mi riferirò usando il simbolo  $\bar{1}$ .

La cifra  $\bar{1}$  funziona esattamente come le altre cifre, soltanto che il suo contributo al valore del numerale è negativo. Questo quindi amplia l'insieme dei numeri rappresentabili tramite stream non più all'intervallo  $[0, 1]$  ma all'intervallo  $[-1, 1]$ .

L'uso di questa cifra comporta anche l'introduzione di una ridondanza in questa notazione. Ridondanza significa che alcuni numeri possono essere scritti non più in un unico modo, ma in molti modi diversi.

In questa notazione, infatti, gli unici due numeri che mantengono una rappresentazione unica sono:

$$1 = \left[ \vec{1} \right] \quad \text{e} \quad \bar{1} = \left[ \vec{\bar{1}} \right]$$

Ad esempio il numero 0 può essere espresso nei seguenti modi (e non solo):

$$0 = \left[ \vec{0} \right] \quad \text{oppure} \quad 0 = \left[ 1 : \vec{\bar{1}} \right] \quad \text{oppure} \quad 0 = \left[ \bar{1} : \vec{1} \right]$$

Infatti:

$$\left[ 1 : \vec{\bar{1}} \right] = \frac{1}{2} + \sum_{i=2}^{\infty} \bar{1} \times 2^{-i} = \frac{1}{2} - \frac{1}{4} - \frac{1}{8} - \frac{1}{16} - \dots = 0$$

$$\left[ \bar{1} : \vec{1} \right] = -\frac{1}{2} + \sum_{i=2}^{\infty} 1 \times 2^{-i} = -\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \dots = 0$$

Va inoltre fatto notare che in questa rappresentazione esistono anche le seguenti identità:

$$1 : \bar{1} : x = 0 : 1 : x$$

$$\bar{1} : 1 : x = 0 : \bar{1} : x$$

Uno dei vantaggi di questa rappresentazione è che permette di eseguire anche quei calcoli che nella rappresentazione classica necessiterebbero di una quantità infinita di input per produrre anche solo una cifra di output. Questo succede perché grazie alla cifra  $\bar{1}$  si possono aggiustare eventuali errori di eccesso avvenuti durante i precedenti passi del calcolo. Questi errori in questa rappresentazione sono comuni in quanto, come detto in precedenza, gli algoritmi sugli stream devono produrre cifre a partire da sinistra e non da destra. Questo porta spesso alla produzione di una cifra, ad esempio 1, che però potrebbe risultare più avanti una cifra troppo grande in quella posizione. Questo errore, che non si potrebbe correggere nella rappresentazione classica, si può invece aggiustare in questa rappresentazione grazie alla cifra  $\bar{1}$ , che permette di diminuire il valore del numero in output.

Tornando al caso della somma tra gli stream

$$1 : 1 : 1 : \dots \quad \text{e} \quad 0 : 0 : 0 : \dots$$

che non poteva essere computata nella rappresentazione classica: il problema stava nel fatto che la prima cifra di questa somma può essere:

- 1 nel caso in cui il valore della somma dovesse risultare  $\geq 1$
- 0 altrimenti

soltanto che per saperlo si sarebbe dovuta analizzare una parte potenzialmente infinita di input. Questo problema non c'è più a seguito dell'introduzione della cifra  $\bar{1}$  in quanto, anche dando in output 1 come prima cifra, nel caso in cui il risultato reale si rivelasse più piccolo di quello atteso, si possono dare in output una o più cifre  $\bar{1}$  per ridurre il valore del risultato calcolato e farlo arrivare al giusto valore della somma.

Alcune proprietà dei numeri binari vengono mantenute in questa rappresentazione mentre altre vengono migliorate. Ad esempio in questa rappresentazione si possono facilmente rappresentare numeri negativi e si può passare da un numero al suo opposto semplicemente negandone le cifre:

$$\begin{aligned} \llbracket 1 : \bar{1} : 0 : 1 : 0 : \vec{0} \rrbracket &= 0.3125 \\ \llbracket \bar{1} : 1 : 0 : \bar{1} : 0 : \vec{0} \rrbracket &= -0.3125 \end{aligned}$$

Anche le proprietà della divisione e moltiplicazione per 2 tramite shift valgono, con l'accortezza che i numeri in questa rappresentazione possono trovarsi soltanto nell'intervallo  $[-1, 1]$ , e quindi non è sempre detto che si possa moltiplicare per 2 un numero senza causare overflow:

$$\begin{aligned} \llbracket 1 : \bar{1} : 0 : 1 : 0 : \vec{0} \rrbracket &= 0.3125 \\ \llbracket 0 : 1 : \bar{1} : 0 : 1 : \vec{0} \rrbracket &= 0.15625 = \frac{0.3125}{2} \\ \llbracket \bar{1} : 0 : 1 : 0 : 0 : \vec{0} \rrbracket &= -0.375 \neq 0.3125 \cdot 2 = 0.625 \end{aligned}$$

È possibile però notare che, anche se il numero 0.3125 è minore di 0.5, la sua moltiplicazione causa comunque un overflow, nonostante il numero:

$$0.3125 \cdot 2 = 0.625$$

sia rappresentabile nel range  $[-1, 1]$ . Questo succede a causa della ridondanza introdotta dalla cifra  $\bar{1}$ . Possiamo notare infatti che le prime due cifre dello stream:

$$1 : \bar{1} : 0 : 1 : 0 : \vec{0}$$

possono essere cambiate in:

$$0 : 1 : 0 : 1 : 0 : \vec{0}$$

grazie alle due uguaglianze indicate in precedenza. Ora, con questa nuova forma dello stesso numero, ci è possibile eseguire la moltiplicazione per 2 tramite shift ottenendo il risultato corretto:

$$\llbracket 1 : 0 : 1 : 0 : \vec{0} \rrbracket = 0.625$$

Gli algoritmi relativi a questa rappresentazione utilizzano queste uguaglianze per manipolare gli stream in input e risolvere questi errori.

Per quanto riguarda il calcolo di numeri reali sull'intera linea e non soltanto nell'intervallo  $[-1, 1]$ , si usa una rappresentazione simile a mantissa-esponente:

$$\llbracket x \rrbracket = \llbracket (m, e) \rrbracket = \llbracket m \rrbracket \times 2^e$$

dove  $x$  è un numerale,  $m$  è uno stream di cifre binarie con segno ed  $e$  è un numero intero di grandezza arbitraria

### 1.3.4 Basi non naturali

Non è necessario, però, che le basi delle rappresentazioni si limitino ai numeri naturali: è possibile infatti utilizzare numeri razionali o reali come base, gli unici vincoli sono che il numero sia computabile e che sia strettamente minore del numero di cifre utilizzate nella rappresentazione. Grazie a queste due regole è possibile definire una rappresentazione in base reale nel seguente modo:

Dato un numero naturale  $d$  e un numero computabile  $b$  tali che  $1 < b < d$ , la rappresentazione del numero  $x$  con  $d$  cifre in base  $b$  è data dalla sequenza di interi  $z_0 : z_1 : z_2 : \dots$  tali che:

$$\forall i \geq 1 . 0 \leq z_i < d$$

$$x = \sum_{i=0}^{\infty} z_i \times b^{-i}$$

Dove il primo intero della sequenza non è vincolato in grandezza e serve a rappresentare la parte intera del numero, mentre il resto delle cifre rappresentano la parte decimale.

Dal primo vincolo imposto si può dedurre che le rappresentazioni in base reale sono numerabili, ma non tutte sono adatte alla computazione dei reali. Tutte quelle con base razionale, infatti, non assicurano di poter effettuare qualunque computazione con una quantità finita di input. Per ottenere questa proprietà è necessario usare un numero irrazionale come base della rappresentazione [1].

### 1.3.5 Rappresentazione in base $\phi$

Il numero irrazionale scelto come base per la rappresentazione dei reali è  $\phi$ . La scelta di questo numero deriva dal fatto che  $\phi$ , sotto opportune ipotesi, è quello che offre gli algoritmi più semplici per le operazioni aritmetiche [1].

$$\phi = \frac{\sqrt{5} + 1}{2}$$

L'utilizzo del numero  $\phi$  come base rende possibile la rappresentazione e il calcolo dei numeri reali usando solamente due cifre: 0 e 1. I numeri in base  $\phi$ , dunque, appaiono come dei semplici numeri binari, soltanto che il peso di ogni cifra non è una potenza di 2 ma una potenza di  $\phi$ .

Anche questa base, come la rappresentazione con cifre negative, presenta una grande ridondanza. Esiste infatti la seguente identità:

$$1.00 = 0.11 = 0.\overline{10}$$

E questo vale in qualunque posizione si trovino le cifre 1, 0, 0. Ad esempio:

$$\begin{aligned} 10.0 &= 1.1 &= 1.0\overline{10} &= \phi \\ 100.0 &= 11.0 &= 10.\overline{10} &= \phi^2 \\ 0.11 &= 0.1011 &= 0.\overline{10} &= 1 \end{aligned}$$

Questo è vero per via del fatto che:

$$\phi^{n+2} = \phi^{n+1} + \phi^n$$

Da queste uguaglianze deriva che ogni numero con una rappresentazione finita in base  $\phi$  (a parte 0) può essere scritto in forma non finita ma periodica, con parte periodica diversa da 0.

L'utilizzo di questa base irrazionale comporta anche un importante cambiamento nella forma dei numeri rappresentati. Basta pensare, ad esempio, che lo stesso numero  $\phi$ , che in qualunque base naturale non si può rappresentare con un numero finito di cifre, in questa base si rappresenta con il semplice numero 10. Infatti:

$$10_\phi = 1 \cdot \phi^1 + 0 \cdot \phi^0 = \phi$$

Ovviamente anche ogni numero naturale in base  $\phi$  può essere rappresentato con un numero finito di cifre, dal momento che i numeri naturali si possono ottenere sommando 1 un numero arbitrario (e finito) di volte. Ad esempio:

$$\begin{aligned} 1_{10} &= 0.0 &+ 1.0 &= 1.0 &= 0.11 \\ 2_{10} &= 1.0 &+ 0.11 &= 1.11 &= 10.01 \\ 3_{10} &= 10.01 &+ 1.0 &= 11.01 &= 100.01 \\ 4_{10} &= 100.01 &+ 1.0 &= 101.01 &= 101.0011 \\ 5_{10} &= 101.0011 &+ 0.11 &= 101.1111 &= 110.0111 \end{aligned}$$

Per quanto riguarda la rappresentazione dei numeri reali in questa base, anche in questo caso si usano stream infiniti di cifre. Questa volta, però, i numeri rappresentabili tramite questi stream si trovano nell'intervallo  $[0, \phi]$ . In questa rappresentazione gli unici numeri che si possono scrivere in un unico modo sono

$$\phi = \overrightarrow{1} \quad \text{e} \quad 0 = \overrightarrow{0}$$

mentre tutti gli altri sono rappresentabili in infiniti modi grazie all'identità tipica di questa base.

### Numeri di Fibonacci

Una curiosità di questa base è la facilità con cui si possono rappresentare e calcolare i numeri di Fibonacci. In generale, i numeri di Fibonacci sono dati dalla seguente funzione ricorsiva:

$$\begin{cases} Fib(1) = 1 \\ Fib(2) = 1 \\ Fib(n) = Fib(n-1) + Fib(n-2) \end{cases}$$

oppure dalla formula di Binet:

$$Fib_n = \frac{\phi^n}{\sqrt{5}} - \frac{(1-\phi)^n}{\sqrt{5}} = \frac{\phi^n - (-\phi)^{-n}}{\sqrt{5}}$$

Queste due formule funzionano indipendentemente dalla base che si usa, ma la rappresentazione in base  $\phi$  si caratterizza perché in questa base i numeri di Fibonacci seguono un pattern fatto da una cifra 1 seguita da tre cifre 0, e inoltre i numeri che occupano una posizione pari nella sequenza di Fibonacci si distinguono da quelli che occupano una posizione dispari in base al modo in cui finiscono. Il pattern e la differenza tra posizioni pari e dispari sono facilmente individuabili in questa porzione di sequenza di Fibonacci:

1	=	0.11	= 1	(1)
1	=	1	= 1	(2)
10.01	=	10.0011	= 2	(3)
100.01	=	100.01	= 3	(4)
1000.1001	=	1000.100011	= 5	(5)
10001.0001	=	10001.0001	= 8	(6)
100010.001001	=	100010.00100011	= 13	(7)
1000100.010001	=	1000100.010001	= 21	(8)
10001000.10001001	=	10001000.1000100011	= 34	(9)
100010001.00010001	=	100010001.00010001	= 55	(10)
1000100010.0010001001	=	1000100010.001000100011	= 89	(11)
10001000100.0100010001	=	10001000100.0100010001	= 144	(12)

Da questi esempi si può facilmente notare che il passaggio da un numero al successivo viene fatto principalmente con uno shift a sinistra e, in base alla posizione nella sequenza (pari o dispari), con l'aggiunta o la rimozione di alcune cifre a destra del numero.

Di seguito la dimostrazione che il pattern vale per ogni numero di Fibonacci in base  $\phi$ , fatta per induzione e basandosi sul fatto che il pattern è diverso per numeri in posizione pari e dispari.

**Teorema 1.** *Un qualunque numero di Fibonacci  $Fib(n)$  con  $n > 2$  può essere ottenuto dal precedente nel seguente modo:*

*Se  $n$  è dispari allora  $Fib(n)$  si ottiene facendo uno shift a sinistra delle cifre di  $Fib(n-1)$  e aggiungendo le cifre 00011 a destra.*

*Se  $n$  è pari allora  $Fib(n)$  si ottiene facendo uno shift a sinistra delle cifre di  $Fib(n-1)$  e rimuovendo l'ultima cifra 1 a destra.*

*Dimostrazione.* I numeri di Fibonacci  $Fib(n)$  con  $n \leq 12$  soddisfano la proprietà, come mostrato nella tabella in alto. Per i numeri con  $n > 12$  si procede come segue:

Se  $n$  è dispari il numero si deve ottenere tramite shift a sinistra e aggiunta di cifre 00011 a destra.

Supponiamo che i numeri fino a  $Fib(n-1)$  soddisfino il pattern. Allora i numeri  $Fib(n-1)$  e  $Fib(n-2)$  hanno la seguente forma:

$$\begin{array}{rcl} 100010 \dots 00100011 & & (n-2) \\ 1000100 \dots 010001 & & (n-1) \end{array}$$

Dalla cui somma si ottiene il numero:

$$1100110 \dots 01100111 = 10001000 \dots 1000100011$$

Che messo in colonna con i due precedenti mostra come ci sia stato uno shift a sinistra e aggiunta di cifre 00011 a destra rispetto al precedente:

$$\begin{array}{rcl} 100010 \dots 00100011 & & (n-2) \\ 1000100 \dots 010001 & & (n-1) \\ 10001000 \dots 1000100011 & & (n) \end{array}$$

Se  $n$  è pari il numero si deve ottenere tramite shift a sinistra e rimozione dell'ultima cifra 1 a destra.

Supponiamo che i numeri fino a  $Fib(n-1)$  soddisfino il pattern. Allora i numeri  $Fib(n-1)$  e  $Fib(n-2)$  hanno la seguente forma:

$$\begin{array}{rcl} 100010 \dots 0010001 & & (n-2) \\ 1000100 \dots 01000100011 & & (n-1) \end{array}$$

Dalla cui somma otteniamo il numero:

$$1100110 \dots 01100110011 = 10001000 \dots 100010001$$

Che messo in colonna con i due precedenti mostra come ci sia stato uno shift a sinistra e rimozione dell'ultima cifra 1 rispetto al precedente:

$$\begin{array}{rcl} 100010 \dots 0010001 & & (n-2) \\ 1000100 \dots 01000100011 & & (n-1) \\ 10001000 \dots 100010001 & & (n) \end{array}$$

□

Questa proprietà deriva dall'identità dei numeri in base  $\phi$ :

$$1.00 = 0.11$$

che provoca uno slittamento a sinistra di due cifre 1 consecutive. Nel caso dei numeri di Fibonacci, visto che ogni numero è ottenuto dal precedente facendo uno shift a sinistra, ogni somma di due numeri di Fibonacci consecutivi si trova ad avere una coppia di 1 dove nei precedenti c'erano le singole cifre 1, e questo provoca lo shift a sinistra.

Usando questa proprietà si può ottenere facilmente un numero di Fibonacci arbitrario senza eseguire nessuna operazione ricorsiva ma sfruttando solamente la posizione nella sequenza che si vuole calcolare.

In alternativa, volendo sfruttare la ricorsione, si può calcolare un numero di Fibonacci arbitrario senza eseguire nessuna operazione matematica sui risultati ricorsivi, ma usando soltanto la concatenazione di cifre ai risultati ottenuti ricorsivamente.

Entrambe queste funzioni sono state implementate in Haskell.



## Capitolo 2

# La Golden Notation

La rappresentazione in base  $\phi$ , dunque, permette di rappresentare numeri nell'intervallo  $[0, \phi]$ . Questa notazione verrà d'ora in avanti chiamata notazione semplificata, e la sua interpretazione da numerale a numero verrà fatta così [1]:

$$\llbracket \alpha \rrbracket_s = \sum_{i=1}^{\infty} \alpha_i \times \phi^{-i}$$

Volendo espandere la rappresentazione in base  $\phi$  a tutta la linea dei reali, si può usare anche in questo caso una rappresentazione simile a mantissa-esponente, ma c'è un problema: l'utilizzo di stream infiniti in base  $\phi$  non permette di rappresentare i numeri negativi. Per rimediare a questa mancanza si potrebbe pensare di ampliare la rappresentazione mantissa-esponente a una rappresentazione del tipo segno-mantissa-esponente. Questo purtroppo non è possibile in quanto, analogamente al problema della somma

$$3 : 3 : 3 : \dots + 6 : 6 : 6 : \dots$$

per decidere il segno del risultato di un'operazione servirebbe sapere se è minore o maggiore di zero, ma questo non è possibile quando si lavora con stream infiniti.

La soluzione è quindi modificare il modo in cui una coppia mantissa-esponente venga interpretata, aggiungendo un termine che permette di rappresentare numeri negativi. Questa notazione verrà d'ora in avanti chiamata notazione completa, e la sua interpretazione da numerale a numero verrà fatta così:

$$\llbracket z : \alpha \rrbracket_f = (-1 + \llbracket \alpha \rrbracket_s) \times \phi^{2z} = \left( -1 + \sum_{i=1}^{\infty} \alpha_i \times \phi^{-i} \right) \times \phi^{2z}$$

dove  $\alpha$  è uno stream nella notazione semplificata e  $z$  è un intero di grandezza arbitraria. Nella notazione completa un numerale è uno stream il cui primo elemento è l'esponente, ma questa notazione ha un significato del tutto equivalente alla coppia  $(z, \alpha)$ .

Il termine  $-1$  viene usato per separare i valori di uno stream nella notazione semplificata tra positivi e negativi. In questo modo tutti i valori nell'intervallo  $[0, 1]$  diventano negativi mentre i valori nell'intervallo  $(1, \phi]$  rimangono positivi. Il valore 1 nella notazione semplificata diventa 0 nella notazione completa, indipendentemente dall'esponente.

Gli algoritmi sono divisi tra notazione semplificata e notazione completa. Quelli per la notazione semplificata, non essendo chiusi sull'intervallo  $[0, \phi]$ , restituiscono un valore diviso per un coefficiente di  $\phi$  oppure  $\phi^2$  in base a quanti bit di overflow possono essere generati dall'algoritmo. Quelli per la notazione completa fanno uso di quelli per la notazione semplificata e di alcune operazioni sugli esponenti.

Come detto in precedenza, gli algoritmi sugli stream lavorano da sinistra a destra utilizzando un certo numero di cifre di lookahead dalla testa dello stream e fanno uso delle identità tipiche della rappresentazione per modificare il contenuto di uno stream senza cambiarne il valore.

## 2.1 Addizione semplificata

L'algoritmo per la notazione semplificata fa uso di un massimo di due cifre di lookahead e utilizza due cifre di riporto con peso relativamente  $\frac{d_1}{\phi}$  e  $\frac{d_2}{\phi^2}$ . Il suo risultato è diviso per  $\phi^2$  per far rientrare due eventuali cifre di overflow. Per l'addizione nella notazione semplificata vale dunque la seguente uguaglianza:

$$\llbracket A(\alpha, \beta, a, b) \rrbracket_s = \frac{\llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s + \frac{a}{\phi} + \frac{b}{\phi^2}}{\phi^2}$$

Le regole basate sulle cifre di lookahead sono:

$$\begin{aligned} A(0 : \alpha, 0 : \beta, 0, b) &= 0 : A(\alpha, \beta, b, 0) \\ A(0 : 0 : \alpha, 0 : \beta, 1, b) &= 0 : A(b : \alpha, \beta, 1, 1) \\ A(0 : 1 : \alpha, 0 : 1 : \beta, 1, 1) &= 1 : 0 : A(\alpha, \beta, 0, 1) \\ A(0 : 0 : \alpha, 1 : 0 : \beta, 1, 0) &= 0 : 1 : A(\alpha, \beta, 1, 0) \\ A(0 : \alpha, 1 : \beta, 1, 1) &= 1 : A(\alpha, \beta, 0, 0) \\ A(1 : \alpha, 1 : \beta, 1, b) &= 1 : A(\alpha, \beta, b, 1) \\ \\ A(1 : \alpha, 0 : \beta, a, b) &= A(0 : \alpha, 1 : \beta, a, b) \\ A(\alpha, 1 : \beta, 0, b) &= A(\alpha, 0 : \beta, 1, b) \\ A(a_1 : 1 : \alpha, b_1 : 0 : \beta, a, b) &= A(a_1 : 0 : \alpha, b_1 : 1 : \beta, a, b) \\ A(\alpha, b_1 : 1 : \beta, a, 0) &= A(\alpha, b_1 : 0 : \beta, a, 1) \end{aligned}$$

## 2.2 Addizione completa

L'algoritmo per l'addizione completa si basa su quello per l'addizione semplificata e su operazioni sull'esponente, facendo distinzione in base alla relazione che c'è tra gli esponenti degli addendi:

$$A'(z : \alpha, t : \beta) = \begin{cases} (z+1) : A(\alpha, \beta, 1, 0) & \text{se } z = t \\ A'((z+1) : 1 : 0 : \alpha, t : \beta) & \text{se } z < t \\ A'(z : \alpha, (t+1) : 1 : 0 : \beta) & \text{se } t < z \end{cases}$$

L'algoritmo esegue la vera operazione solo nel caso in cui i due esponenti siano uguali. Negli altri casi non fa altro che aumentare il minore dei due esponenti fino a portarli ad una situazione di parità, sfruttando le seguenti uguaglianze:

$$\begin{aligned}
\llbracket (z+1) : 1 : 0 : \alpha \rrbracket_f &= (-1 + \llbracket 1 : 0 : \alpha \rrbracket_s) \times \phi^{2z+2} \\
&= \left( -1 + \frac{1}{\phi} + \frac{\llbracket \alpha \rrbracket_s}{\phi^2} \right) \times \phi^{2z+2} \\
&= \left( -\phi^2 + \phi + \llbracket \alpha \rrbracket_s \right) \times \phi^{2z} \\
&= (-1 + \llbracket \alpha \rrbracket_s) \times \phi^{2z} \\
&= \llbracket z : \alpha \rrbracket_f
\end{aligned}$$

La correttezza del caso in cui gli esponenti sono uguali deriva dalle seguenti uguaglianze:

$$\begin{aligned}
\llbracket A'(z : \alpha, t : \beta) \rrbracket_f &= \llbracket (z+1) : A(\alpha, \beta, 1, 0) \rrbracket_f \\
&= (-1 + \llbracket A(\alpha, \beta, 1, 0) \rrbracket_s) \times \phi^{2z+2} \\
&= \left( -1 + \frac{\llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s + \frac{1}{\phi}}{\phi^2} \right) \times \phi^{2z+2} \\
&= \left( -\phi^2 + \llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s + \frac{1}{\phi} \right) \times \phi^{2z} \\
&= \left( -\phi - 1 + \llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s + \frac{1}{\phi} \right) \times \phi^{2z} \\
&= \left( -1 - \frac{1}{\phi} - 1 + \llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s + \frac{1}{\phi} \right) \times \phi^{2z} \\
&= (-1 + \llbracket \alpha \rrbracket_s) \times \phi^{2z} + (-1 + \llbracket \beta \rrbracket_s) \times \phi^{2z} \\
&= \llbracket z : \alpha \rrbracket_f + \llbracket z : \beta \rrbracket_f
\end{aligned}$$

## 2.3 Complemento semplificato

Per quanto riguarda la sottrazione, non si implementa un vero algoritmo per sottrarre due numeri in notazione semplificata in quanto l'intervallo  $[0, \phi]$  non include i numeri negativi. Si utilizza invece un algoritmo che nega tutte le cifre di uno stream. Il risultato della negazione di uno stream in notazione semplificata è:

$$\llbracket C(\alpha) \rrbracket_s = \phi - \llbracket \alpha \rrbracket_s$$

Questo deriva dal fatto che:

$$\llbracket \alpha \rrbracket_s + \llbracket C(\alpha) \rrbracket_s = \llbracket \vec{1} \rrbracket_s = \phi$$

Dove:

$$\begin{aligned}
C(1 : \alpha) &= 0 : C(\alpha) \\
C(0 : \alpha) &= 1 : C(\alpha)
\end{aligned}$$

## 2.4 Complemento e Sottrazione completi

In questa sezione si illustrano gli algoritmi per complemento e sottrazione nella notazione completa. È importante notare che l'algoritmo per la sottrazione non fa uso dell'algoritmo per il complemento, il quale però verrà usato in seguito per moltiplicazione e divisione.

La funzione per effettuare il complemento di un numero nella notazione completa fa uso di una funzione ausiliaria che si basa sulla funzione  $C$  della notazione semplificata:

$$\begin{aligned} C'(z : \alpha) &= (z + 1) : C_1(\alpha) \\ C_1(0 : \alpha) &= 1 : 1 : 0 : C(\alpha) \\ C_1(1 : 0 : \alpha) &= 1 : 0 : C_1(\alpha) \\ C_1(1 : 1 : \alpha) &= 1 : 0 : 0 : 1 : C(\alpha) \end{aligned}$$

L'algoritmo per la sottrazione completa diretta funziona in modo simile a quello dell'addizione completa e si basa sugli algoritmi per addizione e complemento semplificati:

$$S'(z : \alpha, t : \beta) = \begin{cases} (z + 1) : A(\alpha, C(\beta), 1, 1) & \text{se } z = t \\ S'((z + 1) : 1 : 0 : \alpha, t : \beta) & \text{se } z < t \\ S'(z : \alpha, (t + 1) : 1 : 0 : \beta) & \text{se } t < z \end{cases}$$

Anche questo algoritmo effettua la vera operazione solo nel caso in cui i due esponenti siano uguali mentre negli altri casi non fa altro che incrementare il minore fino a pareggiarli, sfruttando la stessa uguaglianza usata nell'addizione completa. La prova di correttezza nel caso in cui gli esponenti sono uguali è la seguente:

$$\begin{aligned} \llbracket S'(z : \alpha, t : \beta) \rrbracket_f &= \llbracket (z + 1) : A(\alpha, C(\beta), 1, 1) \rrbracket_f \\ &= (-1 + \llbracket A(\alpha, C(\beta), 1, 1) \rrbracket_s) \times \phi^{2z+2} \\ &= \left( -1 + \frac{\llbracket \alpha \rrbracket_s + \llbracket C(\beta) \rrbracket_s + 1}{\phi^2} \right) \times \phi^{2z+2} \\ &= \left( -1 + \frac{\llbracket \alpha \rrbracket_s + \phi - \llbracket \beta \rrbracket_s + 1}{\phi^2} \right) \times \phi^{2z+2} \\ &= \left( -1 + \frac{\llbracket \alpha \rrbracket_s - \llbracket \beta \rrbracket_s + \phi^2}{\phi^2} \right) \times \phi^{2z+2} \\ &= (\llbracket \alpha \rrbracket_s - \llbracket \beta \rrbracket_s) \times \phi^{2z} \\ &= (-1 + \llbracket \alpha \rrbracket_s) \times \phi^{2z} - (-1 + \llbracket \beta \rrbracket_s) \times \phi^{2z} \\ &= \llbracket z : \alpha \rrbracket_f - \llbracket z : \beta \rrbracket_f \end{aligned}$$

## 2.5 Prodotto semplificato

L'algoritmo del prodotto nella notazione semplificata ha un funzionamento simile a quello dell'addizione: sfrutta un massimo di due cifre di lookahead e restituisce

il risultato considerando due eventuali bit di overflow. I casi con le due cifre di lookahead sono quattro: uno per ogni combinazione di due bit, considerando che il primo bit è sempre 1. Al suo interno utilizza l'algoritmo dell'addizione per spezzare il prodotto in una serie di somme:

$$\begin{aligned}
P(0 : \alpha, \beta) &= 0 : P(\alpha, \beta) \\
P(\alpha, 0 : \beta) &= 0 : P(\alpha, \beta) \\
P(1 : 0 : \alpha, 1 : 0 : \beta) &= 0 : A(A(\alpha, \beta, 0, 0), 0 : P(\alpha, \beta), 1, 0) \\
P(1 : 1 : \alpha, 1 : 0 : \beta) &= A(A(0 : \alpha, \beta, 0, 0), 0 : 0 : P(\alpha, \beta), 1, 0) \\
P(1 : 0 : \alpha, 1 : 1 : \beta) &= A(A(\alpha, 0 : \beta, 0, 0), 0 : 0 : P(\alpha, \beta), 1, 0) \\
P(1 : 1 : \alpha, 1 : 1 : \beta) &= A(A(\alpha, \beta, 0, 0), 0 : 0 : P(\alpha, \beta), 1, 1)
\end{aligned}$$

Il risultato del prodotto semplificato è quindi:

$$\llbracket P(\alpha, \beta) \rrbracket_s = \frac{\llbracket \alpha \rrbracket_s \times \llbracket \beta \rrbracket_s}{\phi^2}$$

## 2.6 Prodotto completo

L'algoritmo per il prodotto completo fa uso di quelli per l'addizione per il prodotto semplificato, di entrambi gli algoritmi per il complemento e fa la somma degli esponenti:

$$P'(z : \alpha, t : \beta) = C'((z + t + 2) : A(C(P(\alpha, \beta)), A(\alpha, \beta, 0, 0), 1, 0))$$

Questo algoritmo è diverso da quello in [1] in quanto nella prova di correttezza è presente un errore. La giusta prova di correttezza (da cui deriva l'algoritmo corretto) è la seguente:

$$\begin{aligned}
\llbracket z : \alpha \rrbracket_f \times \llbracket t : \beta \rrbracket_f &= (-1 + \llbracket \alpha \rrbracket_s) \times (-1 + \llbracket \beta \rrbracket_s) \times \phi^{2z+2t} \\
&= (1 + \llbracket \alpha \rrbracket_s \times \llbracket \beta \rrbracket_s - \llbracket \alpha \rrbracket_s - \llbracket \beta \rrbracket_s) \times \phi^{2z+2t} \\
&= ((\phi^3 + \phi + 1) - \phi^3 - \phi + \llbracket \alpha \rrbracket_s \times \llbracket \beta \rrbracket_s - \llbracket \alpha \rrbracket_s - \llbracket \beta \rrbracket_s) \times \phi^{2z+2t} \\
&= -(-\phi^4 + \phi^3 + \phi - \llbracket \alpha \rrbracket_s \times \llbracket \beta \rrbracket_s + \llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s) \times \phi^{2z+2t} \\
&= -(-\phi^4 + \phi^3 - \llbracket \alpha \rrbracket_s \times \llbracket \beta \rrbracket_s + \llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s + \phi) \times \phi^{2z+2t} \\
&= - \left( -1 + \frac{\phi - \frac{\llbracket \alpha \rrbracket_s \times \llbracket \beta \rrbracket_s}{\phi^2} + \frac{\llbracket \alpha \rrbracket_s + \llbracket \beta \rrbracket_s}{\phi^2} + \frac{1}{\phi}}{\phi^2} \right) \times \phi^{2z+2t+4} \\
&= - \left( -1 + \frac{\llbracket C(P(\alpha, \beta)) \rrbracket_s + \llbracket A(\alpha, \beta, 0, 0) \rrbracket_s + \frac{1}{\phi}}{\phi^2} \right) \times \phi^{2z+2t+4} \\
&= -(-1 + \llbracket A(C(P(\alpha, \beta)), A(\alpha, \beta, 0, 0), 1, 0) \rrbracket_s) \times \phi^{2z+2t+4} \\
&= -\llbracket (z + t + 2) : A(C(P(\alpha, \beta)), A(\alpha, \beta, 0, 0), 1, 0) \rrbracket_f \\
&= \llbracket C'((z + t + 2) : A(C(P(\alpha, \beta)), A(\alpha, \beta, 0, 0), 1, 0)) \rrbracket_f \\
&= \llbracket P'(z : \alpha, t : \beta) \rrbracket_f
\end{aligned}$$

## 2.7 Divisione semplificata

L'algoritmo per la divisione semplificata è il più complicato: utilizza fino a tre cifre di lookahead, due funzioni ausiliarie mutuamente ricorsive ed ha una restrizione sul divisore, il quale deve per forza iniziare con un 1 per assicurarsi che il risultato sia abbastanza piccolo da poterlo rappresentare nella notazione semplificata. L'algoritmo si basa sulla divisione Euclidea e ritorna il risultato diviso per  $\phi$ . L'algoritmo è il seguente:

$$\begin{aligned}
D(\alpha, 1 : \beta) &= D_1(0 : 0 : \alpha, C(\beta)) \\
D_1(0 : 0 : \alpha, \beta) &= D_2(A(\alpha, 0 : \beta, 0, 0), 0 : \alpha, \beta) \\
D_1(0 : 1 : \alpha, \beta) &= D_2(A(\alpha, 0 : \beta, 1, 1), 1 : \alpha, \beta) \\
D_1(1 : 0 : \alpha, \beta) &= D_2(A(\alpha, 1 : \beta, 1, 1), \alpha, \beta) \\
D_2(0 : 0 : \gamma, \alpha, \beta) &= 0 : D_1(\alpha, \beta) \\
D_2(0 : 1 : 0 : \gamma, \alpha, \beta) &= 0 : D_1(\alpha, \beta) \\
D_2(0 : 1 : 1 : \gamma, \alpha, \beta) &= 1 : D_1(0 : 0 : \gamma, \beta) \\
D_2(1 : \gamma, \alpha, \beta) &= 1 : D_1(\gamma, \beta)
\end{aligned}$$

Il cui risultato è:

$$\llbracket D(\alpha, 1 : \beta) \rrbracket = \frac{\llbracket \alpha \rrbracket}{\llbracket \beta \rrbracket \times \phi}$$

## 2.8 Divisione completa

Anche l'algoritmo per la divisione completa è tutt'altro che banale: utilizza una funzione ausiliaria, l'algoritmo semplificato, ed entrambi gli algoritmi per il complemento (semplificato e completo). Inoltre è importante sapere che questo algoritmo non termina nel caso in cui si provi a dividere per 0. Le regole sono le seguenti:

$$\begin{aligned}
D'(z : \alpha, t : 0 : \beta) &= D'_1(C'((z - t) : \alpha), 0 : C(\beta)) \\
D'(z : \alpha, t : 1 : 0 : \beta) &= D'(z : \alpha, (t - 1) : \beta) \\
D'(z : \alpha, t : 1 : 1 : \beta) &= D'_1((z - t + 1) : \alpha, \beta) \\
D'_1(z : \alpha, 0 : 0 : \beta) &= D'_1((z + 1) : \alpha, \beta) \\
D'_1(z : \alpha, 0 : 1 : \beta) &= (z + 1) : D_1(A(\alpha, \beta, 0, 0), C(\beta)) \\
D'_1(z : \alpha, 1 : \beta) &= (z + 1) : D_1(A(0 : \alpha, \beta, 0, 1), C(\beta))
\end{aligned}$$

## Capitolo 3

# Implementazione

Volendo mettere in pratica gli algoritmi sopra elencati, il linguaggio funzionale Haskell si presta molto bene allo scopo. A seguire vengono illustrate le parti più importanti dell'implementazione, che nel complesso include i tipi per la rappresentazione e tutti gli algoritmi per la notazione semplificata e quella completa. In aggiunta sono state implementate una funzione per convertire numeri razionali di grandezza arbitraria da base 10 a base  $\phi$  e le due funzioni per i numeri di Fibonacci.

Per convenzione le funzioni che implementano gli algoritmi per la notazione semplificata iniziano con una `s` seguita da una lettera maiuscola.

### 3.1 Il linguaggio Haskell

Il linguaggio scelto per l'implementazione è Haskell, un linguaggio puramente funzionale con delle caratteristiche che si prestano molto allo scopo del calcolo dei numeri reali esatti. Haskell infatti implementa il concetto di lazy evaluation ed è proprio grazie a questo che permette di utilizzare facilmente gli stream infiniti.

La lazy evaluation, o laziness, è la possibilità di non valutare un'espressione fino a quando non è strettamente necessario. Questo permette di configurare una funzione che genera uno stream infinito senza mai generarlo effettivamente tutto. Se ne genera infatti una cifra alla volta, all'occorrenza. Questo permette di calcolare un numero arbitrario di cifre di un numero reale senza doverlo per forza calcolare tutto, ottenendo così una precisione arbitraria su qualunque computazione di numeri reali.

### 3.2 Tipi per la rappresentazione

La rappresentazione dei numeri reali in base  $\phi$  tramite notazione semplificata e completa risulta piuttosto immediata quando si parla di Haskell. Per rappresentare uno stream infinito, infatti, basta semplicemente usare una lista di un certo tipo, nel nostro caso una lista di numeri interi:

```

1  -- un bit dovrebbe essere soltanto 0 e 1 ma,
2  -- non essendoci un modo semplice e immediato per limitare il valore di un intero
3  -- mantenendo intatto l'aspetto del pattern matching,
4  -- ho preferito usare direttamente il tipo Int,
5  -- con la consapevolezza che le uniche cifre valide sono 0 e 1
6  type Bit = Int
7
8  -- per gli stream in notazione semplificata basta usare una lista di bit
9  type SNStream = [Bit]
10
11 -- per la notazione semplificata ho preferito usare una coppia,
12 -- in quanto rispecchia maggiormente la convenzione di mantissa-esponente.
13 -- il tipo Integer rappresenta un intero di grandezza arbitraria:
14 -- e' necessario per non avere limiti nella rappresentazione dell'esponente
15 type FNStream = (Integer, SNStream)

```

### 3.3 Prodotto completo

Metto ora a confronto gli algoritmi per il prodotto in notazione completa nella versione originale sbagliata e in quella corretta:

```

1  -- prodotto in notazione completa
2  -- algoritmo originale sbagliato:
3  -- restituisce valori completamente fuori scala e talvolta di segno opposto
4  multiplication' :: FNStream -> FNStream -> FNStream
5  multiplication' (z, as) (t, bs) = (
6      z + t + 2,
7      sAddition
8          (sMultiplication as bs)
9          (sComplement (sAddition as bs 0 0))
10         1 0
11  )
12
13 -- prodotto in notazione completa
14 -- algoritmo corretto
15 multiplication :: FNStream -> FNStream -> FNStream
16 multiplication (z, as) (t, bs) = complement (
17     z + t + 2,
18     sAddition
19         (sComplement (sMultiplication as bs))
20         (sAddition as bs 0 0)
21     1 0
22 )

```



### 3.4 Costanti

In questa sezione vengono elencate le costanti usate nelle funzioni di conversione. L'algoritmo per ottenere il valore di un numero intero in notazione completa è banale e deriva direttamente dalla definizione di notazione completa (nelle prossime righe il termine “decimale” si riferisce alla rappresentazione con sole cifre a destra dello zero, non alla base decimale):

- (1) si rappresenta il numero in base  $\phi$
- (2) si divide il numero per  $\phi^2$  fino a farlo diventare soltanto decimale, facendolo shiftare a destra di un numero pari di posizioni
- (3) si divide ulteriormente il numero per  $\phi^2$ , facendolo shiftare a destra di altre due posizioni
- (4) si somma il valore 1 nella sua forma decimale: 0.11
- (5) il numero decimale ottenuto è la mantissa del numero in notazione completa
- (6) l'esponente è il numero di divisioni fatte durante il secondo passaggio, più una fatta durante il terzo passaggio

Ecco un esempio dell'algoritmo applicato al numero 6:

$$6 \xrightarrow{1} 111.0111 \xrightarrow{2} 1.110111 \xrightarrow{2} 0.01110111 \xrightarrow{3} 0.0001110111 \xrightarrow{4} 0.1101110111 \xrightarrow{5, 6}$$

$$6 = \llbracket (3, 1 : 1 : 0 : 1 : 1 : 1 : 0 : 1 : 1 : 1 : \vec{0}) \rrbracket_f$$

```
1 zeros :: SNStream
2 zeros = 0:zeros
3
4 ones :: SNStream
5 ones = 1:ones
6
7 minusOne :: FNStream
8 minusOne = (0, zeros)
9
10 zero :: FNStream
11 zero = (0, 1:1:zeros)
12
13 one :: FNStream
14 one = (1, 1:1:1:zeros)
15
16 phi :: FNStream
17 phi = (2, 1:1:1:zeros)
18
19 two :: FNStream
20 two = (2, 1:1:0:1:1:zeros)
21
22 three :: FNStream
23 three = (2, 1:1:1:1:0:1:zeros)
24
25 four :: FNStream
26 four = (3, 1:1:0:1:0:1:0:1:zeros)
27
28 five :: FNStream
29 five = (3, 1:1:0:1:0:1:1:1:1:zeros)
30
31 six :: FNStream
32 six = (3, 1:1:0:1:1:1:0:1:1:1:zeros)
33
34 seven :: FNStream
35 seven = (3, 1:1:1:0:1:1:0:0:0:1:zeros)
36
37 eight :: FNStream
38 eight = (3, 1:1:1:1:0:1:0:0:0:1:zeros)
39
40 nine :: FNStream
41 nine = (3, 1:1:1:1:0:1:1:1:0:1:zeros)
42
43 ten :: FNStream
44 ten = (3, 1:1:1:1:1:1:0:1:0:1:zeros)
```

## 3.5 Conversione

Per effettuare la conversione di un numero da base 10 a notazione completa ho scelto un algoritmo semplice: si basa sul valore posizionale delle cifre e le converte una per una, accumulandone le somme dei prodotti per la potenza di 10 corrispondente. Le funzioni sono due:

la prima converte un singolo numero intero da base 10 a notazione completa;

la seconda prende in input due numeri interi, rappresentanti rispettivamente la parte intera e la parte decimale di un numero razionale, e li converte separatamente. Poi divide per 10 la parte decimale di tante volte quante sono necessarie per farle assumere il suo valore decimale inteso nel numero razionale e infine somma le due parti.

Entrambe le funzioni non hanno limite di grandezza sull'input grazie al tipo numerico `Integer`

### 3.5.1 Conversione di un intero

Per convertire un singolo numero intero da base 10 a notazione completa si converte singolarmente ogni cifra del numero, da destra a sinistra, usando modulo e divisione per 10. Ogni cifra viene moltiplicata per la potenza di 10 che corrisponde alla sua posizione nel numero e si somma alle altre, ottenendo così il risultato finale del numero in notazione completa.

Si sfrutta dunque la seguente proprietà:

$$x_{10} := d_n \dots d_1 d_0$$

$$x_{10} = d_n \cdot 10^n + \dots + d_1 \cdot 10 + d_0$$

$$x_\phi = d_{n\phi} \cdot 10_\phi^n + \dots + d_{1\phi} \cdot 10_\phi + d_{0\phi}$$

per implementare la seguente funzione:

```

1 integerToGolden :: Integer -> FStream
2 integerToGolden x = if signum x /= -1 then real else mul minusOne real where
3   integerToGolden' 0 r c = r
4   integerToGolden' x r c = case (mod x 10) of
5     0 -> integerToGolden' (div x 10)      r              (mul c ten)
6     1 -> integerToGolden' (div x 10) (add r (mul one  c)) (mul c ten)
7     2 -> integerToGolden' (div x 10) (add r (mul two  c)) (mul c ten)
8     3 -> integerToGolden' (div x 10) (add r (mul three c)) (mul c ten)
9     4 -> integerToGolden' (div x 10) (add r (mul four  c)) (mul c ten)
10    5 -> integerToGolden' (div x 10) (add r (mul five  c)) (mul c ten)
11    6 -> integerToGolden' (div x 10) (add r (mul six   c)) (mul c ten)
12    7 -> integerToGolden' (div x 10) (add r (mul seven c)) (mul c ten)
13    8 -> integerToGolden' (div x 10) (add r (mul eight c)) (mul c ten)
14    9 -> integerToGolden' (div x 10) (add r (mul nine  c)) (mul c ten)
15    real = integerToGolden' (abs x) zero one
16    -- queste righe servono soltanto ad accorciare le righe precedenti
17    add = addition
18    mul = multiplication

```

In questa funzione il segno viene inizialmente ignorato per convertire il valore assoluto dell'intero in input. Poi viene valutato il segno e si decide se moltiplicare il risultato per -1 o lasciarlo così.

### 3.5.2 Conversione di un razionale

La conversione di un razionale si basa sulla conversione di due singoli interi: la parte intera e la parte decimale del numero razionale. Queste due parti vengono convertite separatamente e poi sommate, dopo che la parte decimale è stata scalata alla posizione originale dividendola per una potenza di 10.

Si sfrutta dunque la seguente proprietà:

$$x_{10} := d_n \dots d_1 d_0 \cdot d_{-1} d_{-2} \dots d_{-m}$$

$$x_{10} = d_n \cdot 10^n + \dots + d_1 \cdot 10 + d_0 + d_{-1} \cdot 10^{-1} + d_{-2} \cdot 10^{-2} + \dots + d_{-m} \cdot 10^{-m}$$

$$x_{10} = d_n \cdot 10^n + \dots + d_1 \cdot 10 + d_0 + \frac{d_{-1} \cdot 10^{m-1} + d_{-2} \cdot 10^{m-2} + \dots + d_{-m}}{10^m}$$

$$x_\phi = d_{n\phi} \cdot 10_\phi^n + \dots + d_{1\phi} \cdot 10_\phi + d_{0\phi} + \frac{d_{-1\phi} \cdot 10_\phi^{m-1} + d_{-2\phi} \cdot 10_\phi^{m-2} + \dots + d_{-m\phi}}{10_\phi^m}$$

per implementare la seguente funzione:

```

1 | rationalToGolden :: Integer -> Integer -> FNStream
2 | rationalToGolden i d = if sign then real else mul minusOne real where
3 |     digits = (length . show . abs) d
4 |     integralPart = iToG (abs i)
5 |     decimalPart = division (iToG (abs d)) (iterate (mul ten) one !! digits)
6 |     real = addition integralPart decimalPart
7 |     -- queste righe servono soltanto ad accorciare le righe precedenti
8 |     sign = fixedSign i == fixedSign d
9 |     fixedSign x = if signum x == -1 then -1 else 1
10 |     iToG = integerToGolden
11 |     mul = multiplication

```

In questa funzione il segno dei due interi in input viene inizialmente ignorato per ottenere la conversione dei loro valori assoluti. Poi si procede a controllare se i segni dei due numeri siano concordi o discordi: se sono concordi viene restituito un numero positivo, altrimenti negativo. Questo significa che se in input vengono dati due numeri negativi il risultato sarà comunque positivo.

## 3.6 Fibonacci

Come detto in precedenza, grazie alla base  $\phi$  è possibile rappresentare e calcolare facilmente i numeri di Fibonacci. Per calcolarli, infatti, esistono due modi particolari che ci permettono di trovare l' $n$ -esimo numero della sequenza senza ricorsione o, in alternativa, con ricorsione ma senza effettuare nessuna operazione matematica sui numeri precedenti. Questa proprietà ci viene dal fatto che possiamo trattare i numeri in base  $\phi$  come delle liste di cifre e, grazie alle proprietà dei numeri di Fibonacci in questa base, possiamo calcolarli facendo soltanto concatenazioni di cifre in testa ad una lista.

Va notato, però, che per semplicità le due funzioni non restituiscono un numero in notazione completa bensì una coppia  $(e, m)$  dove  $e$  sta per esponente e  $m$  per mantissa: il significato di questa coppia è esattamente quello della notazione esponente-mantissa e una grande differenza rispetto alla notazione completa è che la mantissa in questo caso non è infinita ma contiene soltanto il minimo numero di bit necessari a rappresentare il numero. Volendo trasformare questo risultato in notazione completa basterebbe seguire l'algoritmo nella sezione 3.4.

### 3.6.1 Senza ricorsione

La funzione per calcolare i numeri di Fibonacci senza basarsi sui numeri precedenti (ovvero senza ricorsione) è la seguente:

```

1 fibNoRec :: Integer -> (Integer, [Bit])
2 fibNoRec n = (n - 1, fs) where
3     digits 0 s = s
4     digits x s = 1:0:0:0:digits (x - 1) s
5     -- queste righe servono soltanto ad accorciare le righe precedenti
6     fs = digits (div (n - 1) 2) (if mod n 2 == 0 then [1] else [1, 1])

```

Si basa sul fatto che il numero  $Fib(n)$  ha  $n - 1$  bit nella parte intera e i restanti nella parte decimale, ed è sempre composto dalle cifre 1000 concatenate per  $\frac{n-1}{2}$  volte, come si può notare dalla tabella nella sezione 1.3.5.

Per precisione va detto che in questa funzione la ricorsione viene usata non sulla funzione `fibNoRec` ma sulla funzione `digits`, necessaria soltanto per generare la mantissa.

### 3.6.2 Con ricorsione

La funzione per calcolare i numeri di Fibonacci con ricorsione ma senza nessuna operazione matematica sui risultati precedenti è la seguente:

```

1 fibRec :: Integer -> (Integer, [Bit])
2 fibRec 1 = (0, [1, 1])
3 fibRec 2 = (1, [1])
4 fibRec n = (f + 2, 1:0:0:0:fs) where
5     (f, fs) = fibRec (n - 2)

```

---

Si basa sul fatto che il numero  $Fib(n)$  può essere ottenuto concatenando le cifre 1000 in testa al numero  $Fib(n-2)$  e incrementando di due l'esponente, come si può notare dalla tabella nella sezione 1.3.5.

# Bibliografia

- [1] DI GIANANTONIO, P. A golden ratio notation for the real numbers. (1997).
- [2] HENNEY, K. *97 Things Every Programmer Should Know: Collective Wisdom from the Experts*. O'Reilly Media, 1 edn. (2010). ISBN 9780596809485,0596809484.
- [3] MINSKY, M. L. *Computation: Finite and Infinite Machines*. Prentice Hall (1967). ISBN 9780131655638,0131655639.
- [4] PLUME, D. A calculator for exact real number computation. (1998).