

Grundlagen der Stackorganisation und -adressierung

Das Stack- (Kellerspeicher-) Prinzip ist in der Informatik von grundsätzlicher Bedeutung, namentlich was die Programmiersprachen, die Compiler und die Systemsoftware angeht. Manche Architekturen haben Vorkehrungen, um Stacks zu unterstützen, manche nicht (dann müssen Stacks als normale Datenbereiche vorgesehen werden, deren Verwaltung mit elementaren Befehlen auszuprogrammieren ist). Die Grundprinzipien bleiben stets die gleichen.

Grundlagen

Ein Stack ist eine Speichieranordnung, die eine gewisse Anzahl gleich langer Informationsstrukturen (*Stack-Elemente*) aufnehmen kann. Es gibt keinen wahlfreien Zugriff, sondern die Speichieranordnung wird implizit von einem Adreßzähler (*Stackpointer SP*) adressiert.

Stackzugriffe

Es gibt nur zwei grundlegende Zugriffsabläufe:

- ein *Push*-Ablauf legt ein Element auf den Stack,
- ein *Pop*-Ablauf entnimmt das zuletzt (vom letzten Push) auf den Stack gelegte Element (beim nächsten Pop wird dann das vom vorletzten Push abgelegte Element entnommen usw.).

Die Stack-Organisation wird deshalb gelegentlich auch als LIFO (Last In, First Out) bezeichnet.

Wachstumsrichtung

Es ist eine reine Konventionsfrage, ob bei Push-Abläufen der Inhalt des Stackpointers erhöht und bei Pop-Abläufen vermindert wird oder umgekehrt.

In vielen Architekturen *wachsen Stacks immer in Richtung niederer Adressen*, d. h. der Stackpointer zeigt anfänglich immer auf die höchstwertige Adresse. Sein Inhalt wird bei Push-Abläufen vermindert und bei Pop-Abläufen erhöht.

Adreßzähl- und Zugriffsreihenfolge

Ebenso ist es eine reine Konventionsfrage, ob bei einem Push zunächst der Stackpointer verändert und dann das neue Element gespeichert wird oder umgekehrt. Die typische Auslegung: Der Stackpointer zeigt *immer auf das oberste Element im Stack* (Top of Stack, TOS), nicht auf die erste freie Stackposition. Bei einem Push wird deshalb der Stackpointer-Inhalt zunächst vermindert (Ausdehnungsrichtung!); dann wird das Element gespeichert. Umgekehrt wird bei einem Pop das Element entnommen und dann der Stackpointer-Inhalt erhöht (Zugriffsprinzip: Predecrement/Postincrement).

Stack-relative Adressierung

Es ist oft von Vorteil, wenn man zu Elementen des Stack auch wahlfrei zugreifen kann. So kann man auch untere Elemente im Stack erreichen, ohne die oberen zuvor entfernen zu müssen. Solche Zugriffe beziehen sich zweckmäßigerweise auf den Stackpointer, so daß das erste, zweite usw. Element im Stack für Lese- und Schreibzugriffe zugänglich ist, wobei der Stackpointer nicht verändert wird (explizite Stackzugriffe nach dem Prinzip Basis + Displacement mit dem Stackpointer als Basisadreßregister). Muß beim Atmel ausprogrammiert werden (SP über E-A-Zugriffe holen oder softwareseitiger Stack, z. B. auf Grundlage des Y-Registers).

Variabel lange Stackelemente

In den meisten Architekturen, so sie überhaupt Stacks vorsehen, sind alle Elemente in einem Stack *gleich lang*. Kürzere Angaben werden zwecks Ablage auf dem Stack entsprechend erweitert.

Stack Frames

Ein Stack Frame ist ein fester Bereich im Stack. Er dient vor allem dazu, die statischen Variablen des laufenden Programms aufzunehmen.

Statische und dynamische Variable

Statische Variable werden im Programmtext deklariert (jeder Variablenname wird angegeben, und es wird ihm ein Datentyp zugewiesen). Beispiel (wir verwenden der Anschaulichkeit halber eine an Pascal und Ada orientierte Syntax):

<i>Artikel_Nr</i> : Integer;	-- 4 Bytes (ganze 32-Bit-Binärzahl)
<i>Bezeichnung</i> : String(64);	-- 64 Bytes (Zeichenkette)
<i>Preis</i> : Unpacked_BCD(16);	-- 16 Bytes (BCD-Zahl)
<i>Länge, Breite, Höhe</i> : Small_Integer;	-- je 2 Bytes (ganze 16-Bit-Binärzahlen)
<i>Gewicht, Spezifisches_Gewicht</i> : Float;	-- je 4 Bytes (32-Bit-Gleitkommazahlen)
usw.	

Jeder diese Variablen muß der Compiler entsprechenden Speicherplatz zuweisen.

Dynamische Variable entstehen hingegen im Laufe der Verarbeitung (also ohne daß sie der Programmierer ausdrücklich deklarieren muß). Beispiel: der Programmierer schreibt hin:

```
Gewicht := Länge * Breite * Höhe * Spezifisches_Gewicht;
```

Der Compiler muß diese Formel in eine Folge von Maschinenbefehlen umsetzen (hierbei sind u. a. verschiedene Datentypen ineinander zu wandeln). Da die einzelnen Befehle nur ganz elementare Operationen ausführen können, fallen im Verlauf der Rechnung Zwischenergebnisse an. Dies sind die dynamischen Variablen, die typischerweise auf dem Stack abgelegt werden.

Sowohl statische als auch dynamische Variable werden im Stack untergebracht

Das muß nicht unbedingt so sein, hat sich aber bewährt. Und zwar vor allem deshalb, weil man gern Programme in Programme schachtelt (Unterprogrammtechnik). Dann liegt es nahe, die verfügbare Speicherkapazität im Sinne eines Stack zu verwalten und den Speicher vom oberen Ende her aufzufüllen. Zuerst kommt der Stack Frame des ersten Programms. Darüber (in Richtung zu den niederen Adressen hin) werden die gerade aktuellen dynamischen Variablen auf den Stack gelegt. Wenn nun das Programm ein Unterprogramm aufruft, kommt dessen Stack Frame auf den Stack, darüber werden dessen dynamische Variable abgelegt usw.

Das Zugriffsproblem

Gemäß dem Rechenablauf wächst oder schrumpft der Stack. Andererseits sind aber immer die gleichen statischen Variablen zu adressieren. Würde man sich aber stets auf den Stackpointer (als Basisadresse) beziehen, so würden sich bei jedem Zugriff andere Displacements zu den statischen Variablen ergeben. Deshalb sieht man typischerweise ein weiteres Adreßregister vor, den sog. Frame Pointer oder Base Pointer (Abbildung 1.1).

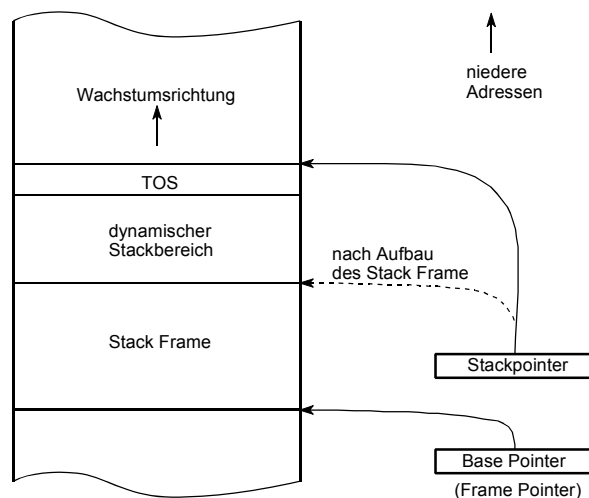


Abbildung 1.1 Stack-Organisation mit Stack Frame

Erklärung zu Abb. 1.1:

Der Base Pointer (Frame Pointer) zeigt stets auf den Anfang des aktuellen Stack Frame (d. h. auf das Wort an der jeweils höchsten Adresse). Alle Inhalte des aktuellen Stack Frame sind somit über negative Displacements (bezogen auf den Base Pointer) erreichbar.

Ruft das aktive Programm seinerseits ein Unterprogramm, so wird der aktuelle Inhalt des Stackpointers in den Base Pointer übernommen, und oberhalb des dynamischen Bereichs des rufenden Programms wird der Stack Frame des gerufenen aufgebaut. Spitzfindigkeiten erläutern wir im folgenden anhand von UNIX .

Grundlagen der systemseitigen Speicherverwaltung

Die Speicherverwaltung hat die Aufgabe, den einzelnen Programmen im Arbeitsspeicher eine angemessene Speicherkapazität zur Verfügung zu stellen. Wieviel Speicher (z. B. in Bytes ausgedrückt) braucht aber ein Programm? - Es sind unterzubringen:

- das Programm selbst,
- die zugehörigen konstanten Daten,
- Arbeits- und Übergabebereiche,
- bedarfsweise Symbol- und Verweistabellen.

In einfachen Systemen kann man die verfügbare Speicherkapazität fest aufteilen (statische Speicheraufteilung). Moderne Hochleistungssysteme sind hingegen dadurch gekennzeichnet, daß sich die Speicherbelegung ständig ändert (dynamische Speicheraufteilung). Abbildung 1.2 veranschaulicht ein Prinzip, das häufig implementiert wird.

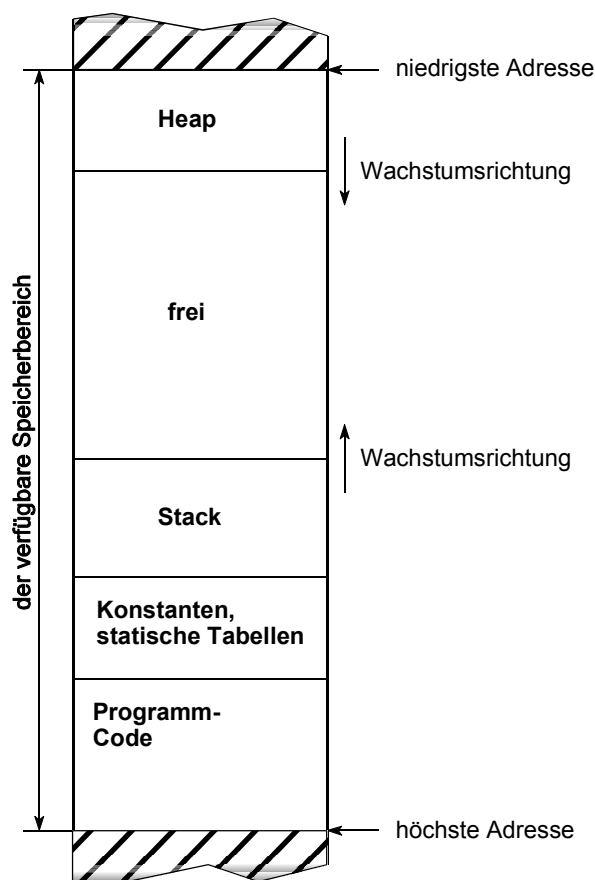


Abbildung 1.2 Zum Prinzip der Speicheraufteilung (Beispiel)

Erklärung zu Abb. 1.2:

Wir beginnen damit, daß zunächst ein "hinreichend" großer Speicherbereich bereitsteht. Dieser wird folgendermaßen belegt.

- der Programmcode an sich wird ganz hinten untergebracht,
- davor kommen die "statischen" - in ihrer Größe unveränderlichen Datenbereiche (Konstanten, Symboltabellen usw.),
- im Anschluß daran - zu den niederen Adressen hin - wird der Stack eingerichtet. Er nimmt dynamische Daten, Parameter, statische Variable, Zwischenergebnisse und Rückkehradressen auf. Er wächst in Richtung niederer Adressen.
- ergänzend zum Stack sieht man oft eine weitere veränderliche Struktur vor, den Heap (sprich: Hiep; wörtlich = Haufen). Der Heap wird am Anfang des Speicherbereichs angeordnet. Er wächst in Richtung höherer Adressen. Zur Verwendung von Stack und Heap siehe Tabelle 1.1.

Sowohl Stack als auch Heap wachsen oder schrumpfen während der Ausführung des Programms. Durch die Anordnung an entgegengesetzten Enden ist stets gewährleistet, daß sich ein möglichst großer freier Bereich zwischen Stack und Heap befindet. Nur in dem - vergleichsweise unwahrscheinlichen - Fall, daß beide Strukturen wachsen und wachsen, kann es vorkommen, daß irgendwann einmal nichts mehr frei ist, daß also der Stack versucht, ein Stück des Heap zu belegen oder umgekehrt. Die Schutzvorkehrungen der Hardware bzw. das Laufzeitsystem der Software sollten dies erkennen und entsprechend reagieren (z. B. mit dem Abbruch der Programmausführung und einer entsprechenden Fehlermeldung). **Das ist aber nicht immer der Fall!!!**

	Stack	Heap
Nutzung (gespeichert werden...)	Rückkehradressen, lokale Daten (verschwinden bei Rückkehr aus der jeweiligen Funktion)	dynamische Daten (bleiben solange erhalten, bis sie explizit (vom Programm) wieder freigegeben werden)
Belegung und Freigabe (Auf- und Abbau)	automatisch gemäß dem LIFO-Prinzip	typischerweise (vgl. Programmiersprache C) vom Programmierer anzufordern und freizugeben
besondere Eignung	für kleinere und einfachere Datenstrukturen (zu beispielsweise 32 oder 64 Bits)	für größere und kompliziertere Datenstrukturen (z. B. von 256 Bytes an aufwärts)

Tabelle 1.1 Zur Verwendung von Stack und Heap

Die UNIX-Stackorganisation

Für jeden Prozeß werden zwei Stacks verwaltet (Abbildung 1.3):

- der User Stack zum Aufrufen von Anwendungsprogrammen,
- der Kernel Stack zum Aufrufen der Systemfunktionen.

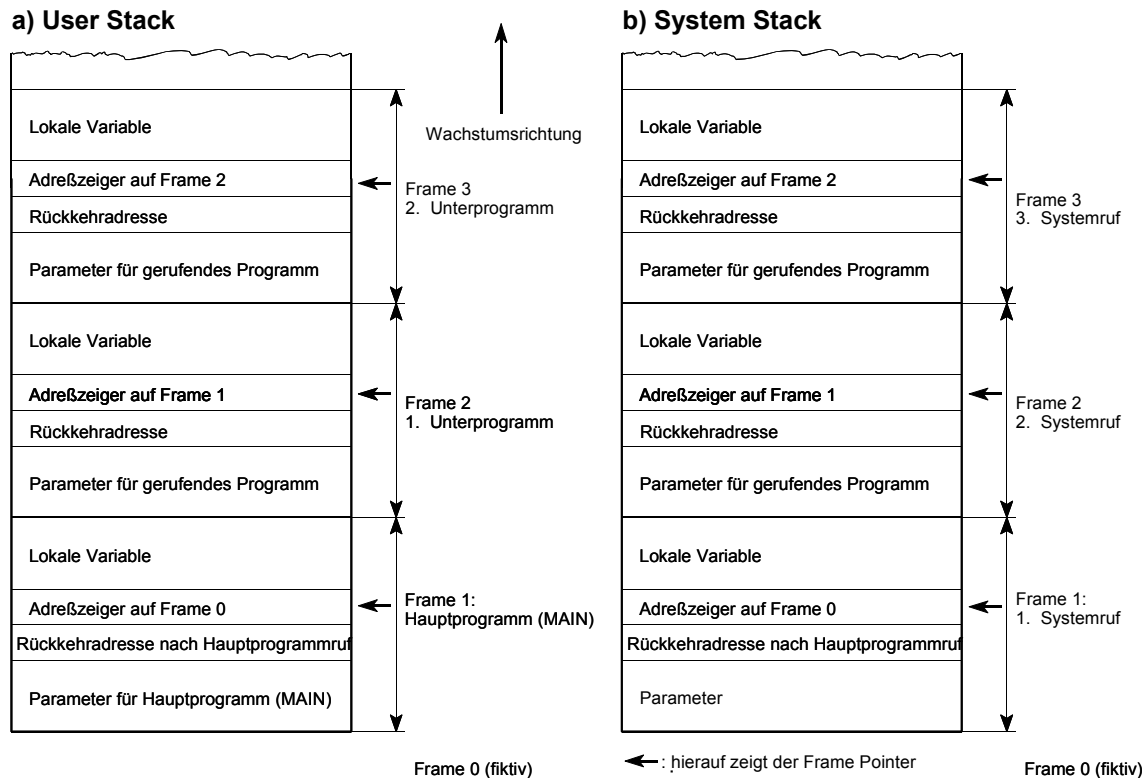


Abbildung 1.3 Die UNIX-Stackorganisation

Erklärung:

Ein UNIX-Programmaufruf läuft folgendermaßen ab:

1. das rufenden Programm legt die zu übergebenden Parameter auf den Stack,
2. der Aufruf wird ausgeführt. Dabei gelangt die Rückkehradresse auf den Stack^{*)}.
3. das gerufene Programm kopiert den bisherigen Frame Pointer auf den Stack (Adreßzeiger als Rückverweis). Typischerweise wird der aktuelle Inhalt des Stackpointers zum neuen Frame Pointer^{**)}.
4. das gerufene Programm kopiert seine lokalen Variablen in den Stack (bzw. schafft auf dem Stack soviel Platz, daß die lokalen Variablen hineinpassen),
5. der aktuelle Frame bzw. Base Pointer wird eingerichtet.

^{*)}: erste Variante: automatisch mittels CALL-Befehl (z. B. IA.32). Zweite Variante: programmseitig, indem der Inhalt des Adreßrettungsregisters auf den Stack gebracht wird (die typischen RISC-Maschinen (Mips, PowerPC, Alpha usw.).

^{**)}: dieser Ablauf wird gelegentlich von der Hardware unterstützt (z. B. IA-32-ENTER-Befehl).

a) rufendes Programm:

PUSH Parameter
 CALL Prozedur (PUSH Rückkehradresse)

b) gerufenes Programm (Funktion, Prozedur)

ENTER-Ablauf (Eintritt):
 PUSH alten Frame Pointer
 Stackpointer wird neuer Frame Pointer (SP => FP)
 DECREMENT SP -- Platz schaffen für lokale Variable

-- der eigentliche Programmablauf --

Rückgabe von Ergebnissen bzw. Funktionswerten: der Parameterbereich ist über den Frame Pointer mit positiven Displacements erreichbar

LEAVE-Ablauf (Rückkehr):
 Stackpointer mit Frame Pointer überladen (FP => SP)
 POP alten Frame Pointer (wird wiederhergestellt)
 RETURN

POP Parameter (Stack säubern)

Abbildung 1.4 Unterprogrammaufruf in einer Laufzeitumgebung, die auf Stack Frames beruht

Typische Konventionen des Unterprogrammrufts

Siehe Tabelle 1.2

	Pascal	C
Reihenfolge der Parameterübergabe	von links nach rechts	von rechts nach links
wer stellt bei der Rückkehr die ursprüngliche Stackbelegung wieder her (Stack Cleanup)?	das gerufene Programm	das rufende Programm
Vor- und Nachteile der Stack-Cleanup-Konvention	<ul style="list-style-type: none"> ▪ Cleanup-Ablauf nur einmal vorhanden (im gerufenen Programm), ▪ Funktionsaufrufe typischerweise nur mit fester Parameteranzahl 	<ul style="list-style-type: none"> ▪ Cleanup-Ablauf in jedem rufenden Programm erforderlich, ▪ erster Parameter (ganz links im Funktionsaufruf) kommt stets auf TOS zu liegen (erleichtert Implementierung von Funktionsaufrufen mit variabler Parameteranzahl)

Tabelle 1.2 Typische Konventionen des Unterprogrammrufts

Die Stackbelegung anhand eines Beispiels

Deklaration einer Funktion:

```
int MAUSI (int A, int B, double C);
```

```
{
int X, Y;
double Z;
float H, I;
...
....
```

```
return (Y);
```

```
}
```

Jetzt wird die Funktion aufgerufen:

```
...
```

```
OMEGA = MAUSI (ALPHA, BETA, GAMMA);
```

```
...
```

Zunächst werden die Parameter auf den Stack gelegt, dann wird die Funktion aufgerufen:

```
PUSH_DOUBLE GAMMA -- Übergabe beginnt von hinten (C-Konvention)
PUSH BETA
PUSH ALPHA
CALL MAUSI
```

Eintritt in die Funktion. Es wird zunächst der Eintrittscode ausgeführt (ENTER-Ablauf):

```
MAUSI: PUSH FP      -- Frame Pointer auf Stack
      MOV  SP, FP   -- neuen FP einrichten
      DEC  SP, 24   -- die lokalen Variablen brauchen 24 Bytes
```

.... jetzt kommt der Funktionskörper von MAUSI

Rückkehr (LEAVE-Ablauf):

```
MOV Y, (FP + 20) -- Rückgabewert in Stack schreiben
MOV FP, SP      -- Zurückstellen des Stackpointers
POP FP          -- alten FP aus Stack zurückholen
RETURN
```

Weiter mit dem rufenden Programm:

```
POP_DOUBLE      -- Stack freimachen
POP
POP OMEGA       -- mit dem letzten POP wird der Rückgabewert zugewiesen
```

Programmoptimierung:

- kurze Funktionen nicht rufen, sondern als Inline-Code einfügen,
- Funktionen ohne Parameter und lokale Variable werden nach einem verkürzten Verfahren aufgerufen. Also: nicht nach der reinen Lehre, sondern auf Leistung programmieren...

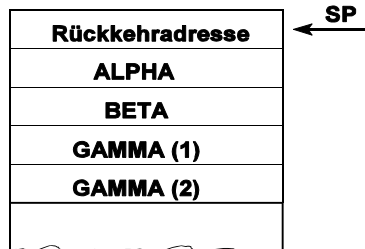
Zum Fehlersuchen (Debugging) müssen diese Optimierungen ggf. ausgeschaltet werden (damit man im Speicherausdruck (Memory Dump) erkennen kann, was eigentlich losgewesen ist):

- kein Inline-Code,
- es wird stets ein richtiger Stack Frame erzeugt.

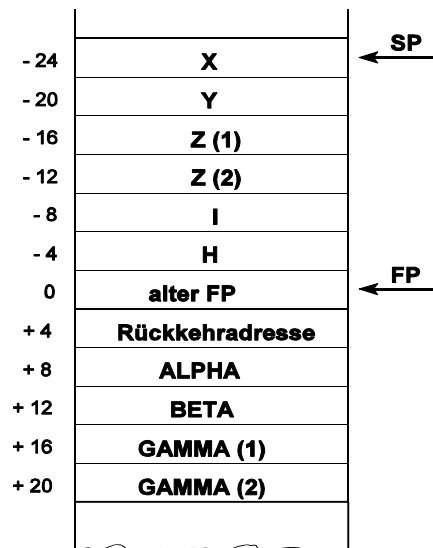
a) Ausgangszustand



b) Stack bei Verzweigung zu MAUSI



c) mit dieser Stackbelegung beginnt die Ausführung des Programmkörpers von MAUSI



Beispiel der Steuerung eines C-Compilers:

