



LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

FAKULTÄT FÜR MATHEMATIK, INFORMATIK UND STATISTIK
INSTITUT FÜR INFORMATIK



Skriptum zur Vorlesung Algorithmische Bioinformatik I&II

gehalten im Sommersemester 2025

und im Wintersemester 2025/26

am Lehrstuhl für Bioinformatik

Volker Heun



30. April 2025

Version 9.02

Vorwort

Dieses Skript entsteht parallel zu den Vorlesungen *Algorithmische Bioinformatik I und II*, die im Sommersemester 25 sowie im Wintersemester 25/26 an der Ludwig-Maximilians-Universität München für Studenten der Bioinformatik und Informatik im Rahmen des von der Ludwig-Maximilians-Universität und der Technischen Universität München gemeinsam veranstalteten Studiengangs Bioinformatik gehalten werden. Dieses Skript basiert dabei in Teilen auf einem früheren Skript zu den Vorlesungen *Algorithmische Bioinformatik I und II*, die im Wintersemester 01/02 und im Sommersemester 02, im Sommer- und Wintersemester 05/06, im Sommer- und Wintersemester 08/09, im Sommer- und Wintersemester 10/11, im Sommer- und Wintersemester 12/13, im Sommer- und Wintersemester 14/15, im Sommer- und Wintersemester 16/17 im Sommer- und Wintersemester 19/20, sowie im Sommer- und Wintersemester 22/23 für Studenten der Bioinformatik und Informatik sowie anderer Fachrichtungen im Rahmen des von der Ludwig-Maximilians-Universität und der Technischen Universität gemeinsam veranstalteten Studiengangs Bioinformatik gehalten wurde.

Teile, die im Sommersemester 25 und Wintersemester 25/26 nicht Teil der Vorlesung waren, sind mit einem Stern (*), Teile, die nur skizziert wurden, sind mit einem (+) markiert.

An dieser Stelle möchte ich mich für die Mithilfe bei der Erstellung des Skript aus dem Jahre 2002 bei folgenden Personen bedanken: Hamed Behrouzi, Michael Engelhardt, Jens Ernst, Peter Lücke, Moritz Maaß, Ingo Rohloff, Sabine Spreer, Hanjo Täubig. Weiterhin möchte ich für die Überarbeitung dieses Skripts insbesondere (in alphabetischer Reihenfolge) Benjamin Albrecht, Constantin Ammar, Joel Daon, Florian Erhard, Johannes Fischer, Caroline Friedel, Simon W. Ginzinger, Markus Gruber, Markus Joppich, Valérie Marot Maximilian Miller, Henrik Otterstedt, Tobias Petri, Korbinian Pürckhauer, Jens Quedenfeld, Katharina Reinisch, Bernhard Schauburger, Konrad Schreiber, Lena Straßer, Jack Wörheide und Simone Wolf für ihre Unterstützung bei den Durchführungen der Veranstaltungen und für Hinweise zu Tippfehlern danken, die somit das aktuell vorliegende Skript erst möglich gemacht haben. Für weitere Hinweise zu Tippfehlern und inhaltlichen Fehlern danke ich Peter Heinig, Anton Smirnov, Thomas Speer, Hanjo Täubig und Stefan Wentzig.

Falls sich dennoch weitere (Tipp)Fehler unserer Aufmerksamkeit entzogen haben sollten, so bin ich für jeden Hinweis darauf (an Volker.Heun@bio.ifi.lmu.de) dankbar.

München, im Sommer- und Wintersemester 2025/26

Volker Heun

Inhaltsverzeichnis

0	Molekularbiologische Grundlagen (*)	1
0.1	Mendelsche Genetik	1
0.1.1	Mendelsche Experimente	1
0.1.2	Modellbildung	2
0.1.3	Mendelsche Gesetze	4
0.1.4	Wo und wie sind die Erbinformationen gespeichert?	4
0.2	Chemische Grundlagen	4
0.2.1	Kovalente Bindungen	5
0.2.2	Ionische Bindungen	7
0.2.3	Wasserstoffbrücken	8
0.2.4	Van der Waals-Kräfte	9
0.2.5	Hydrophobe Kräfte	10
0.2.6	Funktionelle Gruppen	10
0.2.7	Stereochemie und Enantiomerie	11
0.2.8	Tautomerien	13
0.3	DNS und RNS	14
0.3.1	Zucker	14
0.3.2	Basen	16
0.3.3	Polymerisation	18
0.3.4	Komplementarität der Basen	18
0.3.5	Doppelhelix	20
0.4	Proteine	22
0.4.1	Aminosäuren	22

0.4.2	Peptidbindungen	23
0.4.3	Proteinstrukturen	26
0.5	Der genetische Informationsfluss	29
0.5.1	Replikation	29
0.5.2	Transkription	30
0.5.3	Translation	31
0.5.4	Das zentrale Dogma	34
0.5.5	Promotoren	34
0.6	Biotechnologie	35
0.6.1	Hybridisierung	35
0.6.2	Klonierung	36
0.6.3	Polymerasekettenreaktion	36
0.6.4	Restriktionsenzyme	37
0.6.5	Sequenzierung kurzer DNS-Stücke	38
0.6.6	Sequenzierung eines Genoms	40
1	Algorithmik	43
1.1	Einführendes Beispiel: MSS	43
1.1.1	Maximal Scoring Subsequence	43
1.1.2	Naive Lösung	45
1.1.3	Rekursion	48
1.1.4	Dynamische Programmierung	49
1.1.5	Divide-and-Conquer-Ansatz	51
A	Literaturhinweise	59
A.1	Lehrbücher zur Vorlesung	59
A.2	Lehrbücher zur Bioinformatik	60
A.3	Lehrbücher zur Algorithmik und Komplexität	60
A.4	Lehrbücher zur Algorithmenanalyse	61

0.1 Mendelsche Genetik

Dieses Kapitel ist nicht Bestandteil der Vorlesung, sondern dient nur als kurze Einführung in die Molekularbiologie für Hörer ohne besonderen biologischen Hintergrund.

In diesem Einführungskapitel wollen wir uns mit den molekularbiologischen Details beschäftigen, die für die informatische und mathematische Modellierung im Folgenden hilfreich sind. Zu Beginn stellen wir noch einmal kurz die Anfänge der systematischen Genetik, die Mendelsche Genetik, dar.

0.1.1 Mendelsche Experimente

Eine der ersten systematischen Arbeiten zur Vererbungslehre wurde im 19. Jahrhundert von Gregor Mendel geleistet. Unter anderem untersuchte Mendel die Vererbung einer Eigenschaft von Erbsen, nämlich ob die Erbsen eine glatte oder runzlige Oberfläche besitzen. Wie bei allen Pflanzen besitzt dabei jedes Individuum zwei Eltern (im Gegensatz beispielsweise zu Einzellern, die sich durch Zellteilung fortpflanzen).

Bei einer Untersuchung wurden in der so genannten *Elterngeneration* oder *Parental-generation* Erbsen mit *glatter* und Erbsen mit *runzlicher* Oberfläche gekreuzt. Somit hatte in der nachfolgenden Generation, der so genannten *ersten Tochtergeneration* oder *ersten Filialgeneration* jede Erbse je ein Elternteil mit glatter und je ein Elternteil mit runzlicher Oberfläche.

Überraschenderweise gab es bei den Nachkommen der Erbsen in der ersten Tochtergeneration nur noch glatte Erbsen. Man hätte wohl vermutet, dass sowohl glatte als auch runzlige Erbsen vorkommen oder aber leicht runzlige bzw. unterschiedlich runzlige Erbsen auftauchen würden.

Noch überraschender waren die Ergebnisse bei der nachfolgenden Tochtergeneration, der so genannten *zweiten Tochtergeneration* oder *zweiten Filialgeneration*, bei der nun beide Elternteile aus der ersten Tochtergeneration stammten. Hier kamen sowohl glatte als auch wieder runzlige Erbsen zum Vorschein. Interessanterweise waren jedoch die glatten Erbsen im Übergewicht, und zwar im Verhältnis 3 zu 1. Die Frage, die Mendel damals untersuchte, war, wie sich dieses Phänomen erklären lassen konnte.

0.1.2 Modellbildung

Als Modell schlug Gregor Mendel vor, dass die Erbsen für ein bestimmtes Merkmal oder eine bestimmte Ausprägung von beiden Elternteilen je eine Erbinformation erhielt. Im Folgenden wollen wir eine kleinste Erbinformation als *Gen* bezeichnen. Zur Formalisierung bezeichnen wir das Gen, das die glatte Oberfläche hervorruft mit G und dasjenige für die runzlige Oberfläche mit g . Da nun nach unserem Modell jede Erbsen von beiden Elternteilen ein Gen erhält, muss jedes Gen für ein Merkmal doppelt vorliegen. Zwei Erbinformationen, also Gene, die für dieselbe Ausprägung verantwortlich sind, werden als *Allel* bezeichnet. Wir nehmen also an, dass unsere glatten Erbsen in der Elterngeneration die Allele GG und die runzlichen die Allele gg enthalten.

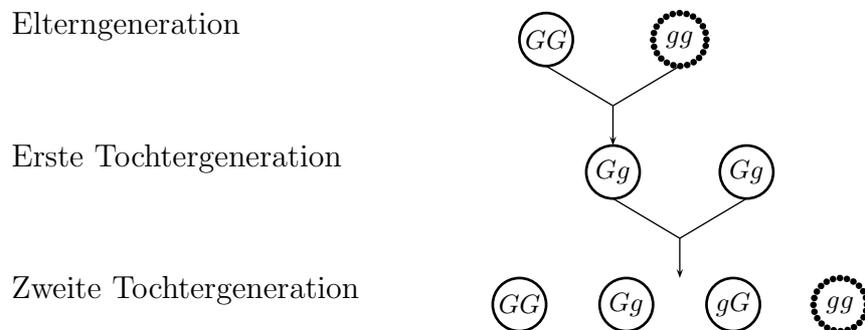


Abbildung 1: Skizze: Mendelsche Vererbung

Welche Erbinformation besitzt nun die erste Tochtergeneration? Sie erhält jeweils ein G und ein g von ihren Eltern und trägt als Erbinformation bezüglich der Oberfläche ein Gg . Was soll nun Gg eigentlich sein? Wir wissen nur, dass GG glatt und gg runzlig bedeutet. Ein Organismus, der bezüglich einer Ausprägung, dieselbe Erbinformation trägt, wird als *reinerbig* oder *homozygot* bezeichnet.

Wir haben nun mit Gg eine *mischerbige* oder *heterozygote* Erbinformation vorliegen. Wie oben bereits angedeutet, könnte die Ausprägung nun gemischt vorliegen, also ein „wenig runzlig“, oder aber eine der beiden Allelen könnte zufällig die Ausprägung bestimmen.

Werden die Merkmale in Mischformen vererbt, wie in „ein wenig runzlig“, dann sagt man, dass das Merkmal *intermediär* vererbt wird. Mitglieder der ersten Tochtergeneration tragen dann also eine Mischung von beidem. Beispielsweise können die Nachfahren von Blumen mit roten bzw. weißen Blüten rosa-farbene Blüten besitzen oder aber auch weiße Blüten mit roten Tupfen etc.

Dies ist aber hier, wie die Experimente von Gregor Mendel gezeigt haben, nicht der Fall: Alle Erbsen der ersten Tochtergeneration sind glatt. Das bedeutet, dass

beide Gene eines Allels gegeneinander konkurrieren und in Abhängigkeit der Gene sich immer eins der beiden als dominant behauptet und den Wettkampf gewinnt. In unserem Falle, setzt sich also das Gen für die glatte Oberfläche gegenüber dem Gen für die runzlige durch. Das Gen, das sich durchsetzt, wird als *dominant* bezeichnet, und dasjenige, das unterliegt, wird als *rezessiv* bezeichnet.

Da nun sowohl die Erbinformation GG als auch Gg für glatte Erbsen stehen, muss man zwischen den so genannten Phänotypen und den Genotypen unterscheiden. Als *Phänotyp* bezeichnet man die sichtbare Ausprägung, also z.B. glatt. Als *Genotyp* bezeichnet man die Zusammensetzung der Erbinformation, also z.B. GG oder Gg für glatte Erbsen. Insbesondere kann also der Genotyp unterschiedlich, aber der Phänotyp gleich sein, wie bei den glatten Erbsen in der Elterngeneration und in der ersten Tochtergeneration.

Wie kann man jetzt die Erscheinung in der zweiten Tochtergeneration erklären? Betrachten wir nun die Eltern, also die Erbsen der ersten Tochtergeneration, die als Genotyp Gg tragen. Nimmt man nun an, dass jedes Elternteil eines seiner Gene eines Allels zufällig (mit gleich hoher Wahrscheinlichkeit) an seine Kinder weitergibt, dann gibt es für die Erbsen der zweiten Tochtergeneration $2 \cdot 2 = 4$ Möglichkeiten, wie sich diese Gene dieses Alles vererben können (siehe Abbildung 2).

	G	g
G	GG	Gg
g	gG	gg

Abbildung 2: Skizze: Vererbung des Genotyps von zwei mischerbigen Eltern

Also sind drei der vier Kombinationen, die im Genotyp möglich sind (GG , gG sowie Gg), im Phänotyp gleich, nämlich glatt. Nur eine der Kombinationen im Genotyp liefert im Phänotyp eine runzlige Erbse. Dies bestätigt in eindrucksvoller Weise das experimentell ermittelte Ergebnis, dass etwa dreimal so viele glatte wie runzlige Erbsen zu beobachten sind.

An dieser Stelle müssen wir noch anmerken, dass diese Versuche nur möglich sind, wenn man in der Elterngeneration wirklich reinerbige Erbsen zur Verfügung hat und keine mischerbigen. Auf den ersten Blick ist dies nicht einfach, da man ja nur den Phänotyp und nicht den Genotyp einfach ermitteln kann. Durch vielfache Züchtung kann man jedoch die Elternteile identifizieren, die reinerbig sind (nämlich, die runzligen sowie die glatten, deren Kinder und Enkelkinder nicht runzlig sind).

0.1.3 Mendelsche Gesetze

Fassen wir hier noch einmal kurz die drei so genannten Mendelschen Gesetze zusammen, auch wenn wir hier nicht alle bis ins Detail erläutert haben:

- 1) **Uniformitätsregel:** Werden zwei reinerbige Individuen einer Art gekreuzt, die sich in einem einzigen Merkmal unterscheiden, so sind alle Individuen der ersten Tochtergeneration gleich.
- 2) **Spaltungsregel:** Werden zwei Mischlinge der ersten Tochtergeneration miteinander gekreuzt, so spalten sich die Merkmale in der zweiten Tochtergeneration im Verhältnis 1 zu 3 bei dominant-rezessiven Genen und im Verhältnis 1 zu 2 zu 1 bei intermediären Genen auf.
- 3) **Unabhängigkeitsregel** Werden zwei mischerbige Individuen, deren Elterngeneration sich in zwei Merkmalen voneinander unterschieden hat, miteinander gekreuzt, so vererben sich die einzelnen Erbanlagen unabhängig voneinander.

Die Unabhängigkeitsregel gilt in der Regel nur, wenn die Gene auf verschiedenen Chromosomen sitzen bzw. innerhalb eines Chromosoms so weit voneinander entfernt sind, dass eine so genannte *Crossing-Over-Mutation* hinreichend wahrscheinlich ist.

0.1.4 Wo und wie sind die Erbinformationen gespeichert?

Damit haben wir die Grundlagen der Genetik ein wenig kennen gelernt. Es stellt sich jetzt natürlich noch die Frage, wo und wie die Gene gespeichert werden. Dies werden wir in den folgenden Abschnitten erläutern.

Zum Abschluss noch ein paar Notationen. Wie bereits erwähnt, bezeichnen wir ein *Gen* als den Träger einer kleinsten Erbinformation. Alle Gene eines Organismus zusammen bilden das *Genom*. Wie bereits aus der Schule bekannt sein dürfte, ist das Genom auf dem oder den *Chromosom(en)* gespeichert (je nach Spezies).

0.2 Chemische Grundlagen

Bevor wir im Folgenden auf die molekularbiologischen Grundlagen näher eingehen, wiederholen wir noch ein paar elementare Begriffe und Eigenschaften aus der Chemie bzw. speziell aus der organischen und der Biochemie. Die in der Biochemie wichtigsten auftretenden Atome sind Kohlenstoff (C), Sauerstoff (O), Wasserstoff (H), Stickstoff (N), Schwefel (S), Kalzium (Ca), Eisen (Fe), Magnesium (Mg), Kalium (K)

und Phosphor (P). Diese Stoffe lassen sich beispielsweise mit folgendem Merkspruch behalten: **COHNS CaFe Mit großem Kuchen-Paket**. Zunächst einmal wiederholen wir kurz die wichtigsten Grundlagen der chemischen Bindungen.

0.2.1 Kovalente Bindungen

Die in der Biochemie wichtigste Bindungsart ist die *kovalente Bindung*. Hierbei steuern zwei Atome je ein Elektron bei, die dann die beiden Atome mittels einer gemeinsamen Bindungswolke zusammenhalten. Im Folgenden wollen wir den Raum, für den die Aufenthaltswahrscheinlichkeit eines Elektrons bzw. eines Elektronenpaares (nach dem Pauli-Prinzip dann mit verschiedenem Spin) am größten ist, als *Orbital* bezeichnen.

Hierbei sind die Kohlenstoffatome von besonderer Bedeutung, die die organische und Biochemie begründen. Die wichtigste Eigenschaft der Kohlenstoffatome ist, dass sie sowohl Einfach-, als auch Doppel- und Dreifachbindungen untereinander ausbilden können. Das Kohlenstoffatom hat in der äußersten Schale 4 Elektronen. Davon befinden sich im Grundzustand zwei in einem so genannten *s-Orbital* und zwei jeweils in einem so genannten *p-Orbital*.

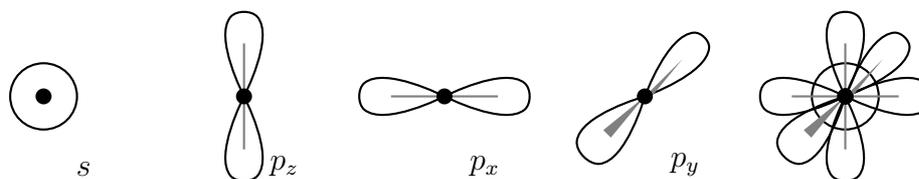


Abbildung 3: Skizze: Räumliche Ausdehnung der Orbitale

Das *s-Orbital* ist dabei kugelförmig, während die drei verschiedenen *p-Orbitale* jeweils eine *Doppelhantel* ausbilden, die paarweise orthogonal zueinander sind. In Abbildung 3 ist die räumliche Ausdehnung des *s-* und der drei *p-Orbitale* schematisch dargestellt, von denen jedes bis zu zwei Elektronen aufnehmen kann. Ganz rechts sind alle Orbitale gleichzeitig zu sehen, die in der Regel für uns interessant sein werden.

In Einfachbindungen befinden sich beim Kohlenstoffatom die einzelnen Elektronen in so genannten *sp³-hybridisierten Orbitalen*, die auch als *q-Orbitale* bezeichnet werden. Hierbei bilden sich aus den 3 Hanteln und der Kugel vier energetisch äquivalente keulenartige Orbitale. Dies ist in der Abbildung 4 links dargestellt. Die Endpunkte der vier Keulen bilden dabei ein Tetraeder aus. Bei einer Einfachbindung überlappen sich zwei der Keulen, wie in Abbildung 4 rechts dargestellt. Die in der Einfachbindung überlappenden *q-Orbitale* bilden dann ein so genanntes *σ-Orbital*.



Abbildung 4: Skizze: sp^3 hybridisierte Orbitale sowie eine Einfachbindung

In Doppelbindungen sind nur zwei p -Orbitale und ein s -Orbital zu drei Keulen hybridisiert, so genannte sp^2 -Orbitale. Ein p -Orbital bleibt dabei bestehen. Dies ist in Abbildung 5 links illustriert. Die in Doppelbindungen überlappenden p -Orbitale werden dann auch als π -Orbital bezeichnet. Bei einer Doppelbindung überlappen sich zusätzlich zu den zwei keulenförmigen hybridisierten q -Orbitalen, die die σ -Bindung bilden, auch noch die beiden Doppelhanteln der p -Orbitale, die dann das π -Orbital bilden. Dies ist schematisch in der Abbildung 5 rechts dargestellt (die Abbildungen sind nicht maßstabsgetreu).

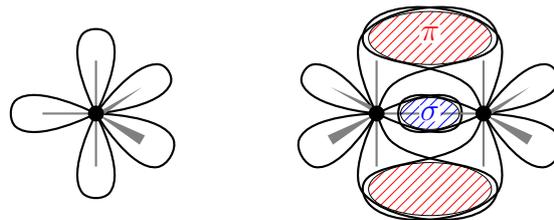


Abbildung 5: Skizze: sp^2 hybridisierte Orbitale und eine Doppelbindung

Die Bindung der Doppelbindung, die durch Überlappung von q -Orbitalen entsteht, wird auch σ -Bindung genannt, die Bindung der Doppelbindung, die durch Überlappung von p -Orbitalen entsteht, wird als π -Bindung bezeichnet. Ähnlich verhält es sich bei Dreifachbindungen, wo zwei p -Orbitale verbleiben und ein s - und nur ein p -Orbital zu einem sp -Orbital hybridisieren. Die konkrete Art der Hybridisierung der Orbitale der Kohlenstoffatome eines bestimmten Moleküls ist deshalb so wichtig, weil dadurch die dreidimensionale Struktur des Moleküls festgelegt wird. Die vier sp^3 -Orbitale zeigen in die Ecken eines Tetraeders, die drei sp^2 -Orbitale liegen in einer Ebene, auf der das verbleibende p -Orbital senkrecht steht, die zwei sp -Orbitale schließen einen Winkel von 180° ein und sind damit gerade gestreckt, auf ihnen stehen die verbleibenden beiden p -Orbitale senkrecht.

Bei zwei benachbarten Doppelbindungen (wie im Butadien, $H_2C=CH-HC=CH_2$) verbinden sich in der Regel die beiden benachbarten Orbitale, die die jeweilige π -Bindung zur Doppelbindung machen, um dann quasi eine π -Wolke über alle vier Kohlenstoffatome auszubilden, da dies energetisch günstiger ist. Daher spricht man

bei den Elektronen in dieser verschmolzenen Wolke auch von *delokalisierten π -Elektronen*.

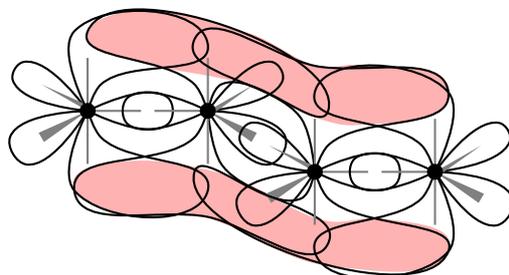


Abbildung 6: Skizze: Delokalisierte π -Bindung im Butadien

Ein Beispiel hierfür ist das *Benzol*-Molekül (C_6H_6). Aus energetischen Gründen bilden sich in dem Ring aus sechs Kohlenstoffatomen nicht drei einzelne alternierende Doppelbindungen aus, sondern eine große Wolke aus sechs delokalisierten π -Elektronen, die die starke Bindung des Benzolrings begründen.

Wie wir später noch sehen werden, kann sich eine Wolke aus delokalisierten π -Elektronen auch aus den π -Elektronen einer $C=C$ Doppelbindung und dem nicht-bindenden Orbital eines Sauerstoff- oder Stickstoffatoms bilden. Stickstoff bzw. Sauerstoff besitzen in der äußersten Schale mehr als vier Elektronen und daher kann sich ein p -Orbital mit zwei Elektronen ausbilden. Dieses hat dann bezüglich der Delokalisation von π -Elektronen ähnliche Eigenschaften wie eine π -Bindung.

Die Energie einer kovalenten Bindung variiert zwischen 200kJ/mol und 450kJ/mol (Kilojoule pro Mol), wobei Kohlenstoffatome untereinander relativ starke Bindungen besitzen (etwa 400kJ/mol). Die Angabe dieser absoluten Werte ist für uns eigentlich nicht von Interesse. Wir geben sie hier nur an, um die Stärken der verschiedenen Bindungsarten im Folgenden vergleichen zu können.

0.2.2 Ionische Bindungen

Bei *ionischen Bindungen* gibt ein Atom, das so genannte *Donatoratom*, ein Elektron an ein anderes Atom, das so genannte *Akzeptoratom*, ab. Damit sind die Donatoratome positiv und die Akzeptoratome negativ geladen. Durch die elektrostatische Anziehungskraft (und die Abstoßung gleichnamiger Ladung) bildet sich in der Regel ein Kristallgitter aus, das dann abwechselnd aus positiv und negativ geladenen Atomen besteht.

Ein bekanntes Beispiel hierfür ist Kochsalz, d.h. Natriumchlorid ($NaCl$). Dabei geben die Natriumatome jeweils das äußerste Elektron ab, das dann von den Chloratomen

aufgenommen wird. Dadurch sind die Natriumatome positiv und die Chloratome negativ geladen, die sich dann innerhalb eines Kristallgitters anziehen.

Hier wollen wir noch deutlich den Unterschied herausstellen, ob wir diese Bindungen in wässriger Lösung oder ohne Lösungsmittel betrachten. Ohne Wasser als Lösungsmittel sind ionische Bindungen sehr stark. In wässriger Lösung sind sie jedoch sehr schwach, sie werden etwa um den Faktor 80 schwächer. Beispielsweise löst sich das doch recht starke Kristallgitter des Kochsalzes im Wasser nahezu auf. In wässriger Lösung beträgt die Energie einer ionischen Bindung etwa 20 kJ/mol.

0.2.3 Wasserstoffbrücken

Eine andere für uns sehr wichtige Anziehungskraft, die keine Bindung im eigentlichen chemischen Sinne ist, sind die *Wasserstoffbrücken*. Diese Anziehungskräfte werden im Wesentlichen durch die unterschiedlichen Elektronegativitäten der einzelnen Atome bedingt.

Die Elektronegativität ist ein Maß dafür, wie stark die Elektronen in der äußersten Schale angezogen werden. Im Periodensystem der Elemente wächst der Elektronegativitätswert innerhalb einer Periode von links nach rechts, weil dabei mit der Anzahl der Protonen auch die Kernladung und damit auch die Anziehungskraft auf jedes einzelne Elektron ansteigt. Innerhalb einer Hauptgruppe sinkt die Elektronegativität mit zunehmender Ordnungszahl, weil die Außenelektronen sich auf immer höheren Energieniveaus befinden und der entgegengesetzt geladene Kern durch die darunterliegenden Elektronen abgeschirmt wird. Deshalb ist z.B. Fluor das Element mit dem größten Elektronegativitätswert.

Eine Liste der für uns wichtigsten Elektronegativitäten in für uns willkürlichen Einheiten ist in Abbildung 7 angegeben. Hier bedeutet ein größerer Wert eine größere Affinität zu Elektronen.

Atom	C	O	H	N	S	P
EN	2,55	3,44	2,20	3,04	2,58	2,19

Abbildung 7: Tabelle: Elektronegativitäten nach Pauling

Bei einer kovalenten Bindung sind die Elektronenwolken in Richtung des Atoms mit der stärkeren Elektronegativität hin verschoben. Dadurch bekommt dieses Atom eine teilweise negative Ladung, während das andere teilweise positiv geladen ist. Ähnlich wie bei der ionischen Bindung, wenn auch bei weitem nicht so stark, wirkt diese Polarisierung der Atome anziehend.

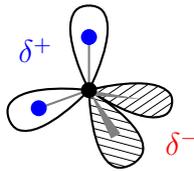


Abbildung 8: Skizze: Polarität bei einem Wassermolekül

Insbesondere Wasser ist für die Ausbildung von zahlreichen Wasserstoffbrücken bekannt. In Abbildung 8 ist ein Wassermolekül schematisch dargestellt.

Wie beim Kohlenstoffatom sind die drei p -Orbitale und das s Orbital zu vier q -Orbitalen hybridisiert. Da das Sauerstoffatom in der äußersten Schale sechs anstatt vier Elektronen besitzt, sind bereits zwei der q -Orbitale des Sauerstoffatoms mit je zwei Elektronen besetzt und können daher keine kovalente Bindung eingehen. Man bezeichnet diese Orbitale daher auch als *nichtbindend*.

Die beiden anderen werden im Wasser gemäß der Formel H_2O mit jeweils einem Wasserstoffatom protoniert. Da nun die beiden nichtbindenden Orbitale (zumindest aus dieser Richtung auf das Sauerstoffatom) negativ geladen sind und die beiden protonierten bindenden Orbitale positiv geladen sind, wirkt das Wasser als Dipol und die Wassermoleküle hängen sich wie viele kleine Stabmagneten aneinander. Einziger Unterschied ist hier dass die Teilladungen in den Ecken eines Tetraeders sitzen, so dass sich die Wassermoleküle ähnlich wie die Kohlenstoffatome im Kristallgitter des Diamanten anordnen.

Im gefrorenen Zustand ist das Kristallgitter von Wasser (also Eis) nahezu ein Diamantengitter, während im flüssigen Zustand die Wasserstoffbrücken häufig aufbrechen und sich wieder neu bilden. Daher ist es zum einen flüssig, und zum anderen kann es im flüssigen Zustand dichter gepackt werden als im gefrorenen Zustand. Erst dadurch nimmt Wasser den flüssigen Zustand bei Zimmertemperatur an, während sowohl Wasserstoff wie auch Sauerstoff einen sehr niedrigen Siedepunkt besitzen. Eine Wasserstoffbrückenbindung ist mit ca. 21 kJ/mol deutlich schwächer als eine kovalente oder eine ionische Bindung.

0.2.4 Van der Waals-Kräfte

Die *Van der Waals-Anziehung* bzw. *Van der Waals-Kräfte* treten insbesondere in großen bzw. langen Molekülen, wie Kettenkohlenwasserstoffen auf. Da der Ort der Elektronen ja nicht festgelegt ist (Heisenbergsche Unschärferelation), können sich durch die Verlagerung der Elektronen kleine Dipolmomente in den einzelnen Bindungen ergeben. Diese beeinflussen sich gegenseitig und durch positive Rückkopplungen

können diese sich verstärken. Somit können sich lange Moleküle fester aneinander legen als kürzere. Dies ist mit ein Grund dafür, dass die homologe Reihe der Alkane (C_nH_{2n+2}) mit wachsender Kohlenstoffanzahl bei Zimmertemperatur ihren Aggregatzustand von gasförmig über flüssig bis zu fest ändert. Die Energie der Van der Waals-Kraft liegt bei etwa 4kJ/mol.

0.2.5 Hydrophobe Kräfte

Nichtpolare Moleküle, wie Fette, können mit Wasser keine Wasserstoffbrücken ausbilden. Dies ist der Grund, warum nichtpolare Stoffe in Wasser unlöslich sind. Solche Stoffe werden auch als *hydrophob* bezeichnet, während polare Stoffe auch als *hydrophil* bezeichnet werden. Aufgrund der Ausbildung von zahlreichen Wasserstoffbrücken innerhalb des Wassers (auch mit hydrophilen Stoffen) tendieren hydrophobe Stoffe dazu, sich möglichst eng zusammenzulagern, um eine möglichst kleine Oberfläche (gleich Trennfläche zum Wasser) auszubilden. Diese Tendenz des Zusammenlagerns hydrophober Stoffe in wässriger Lösung wird als *hydrophobe Kraft* bezeichnet.

0.2.6 Funktionelle Gruppen

Wie schon zu Beginn bemerkt spielt in der organischen und Biochemie das Kohlenstoffatom die zentrale Rolle. Dies liegt insbesondere daran, dass es sich mit sich selbst verbinden kann und sich so eine schier unendliche Menge an verschiedenen Molekülen konstruieren lässt. Dabei sind jedoch auch andere Atome beteiligt, sonst erhalten wir bekanntlich Graphit oder den Diamanten.

Chem. Rest	Gruppe	gewöhnlicher Name
-CH ₃	Methyl	
-OH	Hydroxyl	Alkohol
-NH ₂	Amino	Amine
-NH-	Imino	
-CHO	Carbonyl	Aldehyde
-CO-	Carbonyl	Ketone
-COO-	Ester	Ester
-COOH	Carboxyl	organische Säure
-CN	Cyanid	Nitrile
-SH	Sulfhydril	Thiole

Abbildung 9: Tabelle: Einige funktionelle (organische) Gruppen

Um diese anderen vorkommenden Atome bezüglich ihrer dem Molekül verleihenden Eigenschaften ein wenig besser einordnen zu können, beschreiben wir die am meisten vorkommenden *funktionellen Gruppen*. Die häufigsten in der Biochemie auftretenden einfachen funktionellen Gruppen sind in der Tabelle 9 zusammengefasst.

0.2.7 Stereochemie und Enantiomerie

In diesem Abschnitt wollen wir einen kurzen Einblick in die *Stereochemie*, speziell in die *Enantiomerie* geben. Die Stereochemie beschäftigt sich mit der räumlichen Anordnung der Atome in einem Molekül. Beispielsweise ist entlang einer Einfachbindung die Rotation frei möglich. Bei Doppelbindungen ist diese aufgrund der π -Bindung eingeschränkt und es kann zwei mögliche räumliche Anordnungen desselben Moleküls geben. In Abbildung 10 sind zwei Formen für Äthendiol angegeben. Befinden sich beide (der bedeutendsten) funktionellen Gruppen auf derselben Seite der Doppelbindung, so spricht man vom *cis-Isomer* andernfalls von *trans-Isomer*. Bei der cis-trans-Isomerie kann durch Energiezufuhr die Doppelbindung kurzzeitig

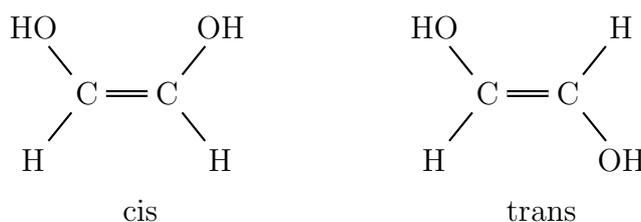


Abbildung 10: Skizze: Cis-Trans-Isomerie bei Äthendiol

geöffnet werden und um 180° gedreht werden, so dass die beiden Isomere ineinander überführt werden können.

Es hat sich herausgestellt, dass scheinbar identische Stoffe (aufgrund der Summen- und Strukturformel) sich unter bestimmten Bedingungen unterschiedlich verhalten können. Dies sind also solche Isomere, die sich nicht ineinander überführen lassen. Betrachten wir dazu in Abbildung 11 ein Kohlenstoffatom (schwarz darge-

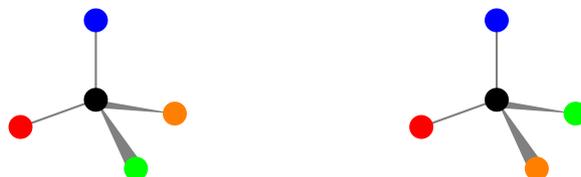


Abbildung 11: Skizze: Asymmetrisches Kohlenstoffatom

stellt) und vier unterschiedliche funktionelle Gruppen (farbig dargestellt), die jeweils

mittels einer Einfachbindung an das Kohlenstoffatom gebunden sind. Das betrachtete Kohlenstoffatom wird hierbei oft als *zentrales Kohlenstoffatom* bezeichnet. Auf den ersten Blick sehen die beiden Moleküle in Abbildung 11 gleich aus. Versucht man jedoch, die beiden Moleküle durch Drehungen im dreidimensionalen Raum zur Deckung zu bringen, so wird man feststellen, dass dies gar nicht geht. Die beiden Moleküle sind nämlich Spiegelbilder voneinander.

Daher werden Moleküle als *chiral* (deutsch Händigkeit) bezeichnet, wenn ein Molekül mit seinem Spiegelbild nicht durch Drehung im dreidimensionalen Raum zur Deckung gebracht werden kann. Die beiden möglichen, zueinander spiegelbildlichen Formen nennen wir *Enantiomere*. Beide Formen werden auch als *enantiomorph* zueinander bezeichnet. Für Kohlenstoffatome, die mit vier *unterschiedlichen* Resten verbunden sind, gilt dies immer. Aus diesem Grund nennt man ein solches Kohlenstoffatom auch ein *asymmetrisches Kohlenstoffatom*.

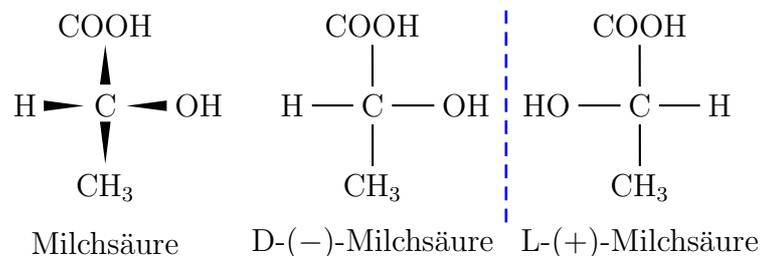


Abbildung 12: Skizze: Milchsäure

Ein einfaches (und bekanntes) Beispiel hierfür ist die Milchsäure. Hierbei sind die funktionellen Gruppen, die an einem zentralen Kohlenstoffatom sitzen, ein Wasserstoffatom, eine Hydroxyl-, eine Methyl und eine Carboxylgruppe. In Abbildung 12 sind rechts die beiden Formeln der beiden spiegelbildlichen Formen dargestellt. In einer zweidimensionalen Abbildung muss man jedoch eine Konvention einführen, wie man die Projektion vornimmt. Die längste Kohlenstoffkette wird dabei von oben nach unten dargestellt. Hier also das zentrale Kohlenstoffatom und das Kohlenstoffatom der Methylgruppe. Dabei wird oben von der charakteristischsten funktionellen Gruppe, hier die Carboxylgruppe, bestimmt. Dabei ist zu verstehen, dass die vertikale Kette hinter dem zentralen Kohlenstoffatom unter der Papierebene verschwindet, während die beiden restlichen Seitenketten aus der Papierebene dem Leser entgegen kommen (dies wird als *Fischer-Projektion* bezeichnet). Dies ist in der Abbildung 12 ganz links noch einmal illustriert.

Da nun im mittleren Teil der Abbildung 12 die bedeutendere funktionelle Gruppe, also die Hydroxylgruppe gegenüber dem Wasserstoffatom, rechts sitzt, wird dieses Enantiomer mit D-Milchsäure (latein. dexter, rechts) bezeichnet. Rechts handelt es sich dann um die L-Milchsäure (latein. laevis, links).

Diese Bezeichnungsweise hat nichts mit den aus der Werbung bekannten links- bzw. rechtsdrehenden Milchsäuren zu tun. Die Namensgebung kommt von der Tatsache, dass eine Lösung von Milchsäure, die nur eines der beiden Enantiomere enthält, polarisiertes Licht dreht. Dies gilt übrigens auch für die meisten Moleküle, die verschiedene Enantiomere besitzen. Je nachdem, ob es polarisiertes Licht nach rechts oder links verdreht, wird es als *rechtsdrehend* oder *linksdrehend* bezeichnet (und im Namen durch (+) bzw. (-) ausgedrückt).

Bei der Milchsäure stimmen zufällig die D- bzw. L-Form mit der rechts- bzw. linksdrehenden Eigenschaft überein. Bei Aminosäuren, die wir noch kennen lernen werden, drehen einige L-Form nach rechts! Hier haben wir einen echten Unterschied gefunden, mit dem sich Enantiomere auf makroskopischer Ebene unterscheiden lassen.

In der Chemie wurde die *DL-Nomenklatur* mittlerweile zugunsten der so genannten *RS-Nomenklatur* aufgegeben. Da jedoch bei Zuckern und Aminosäuren oft noch die DL-Nomenklatur verwendet wird, wurde diese hier angegeben. Für Details bei der DL - sowie der RS-Nomenklatur verweisen wir auf die einschlägige Literatur.

0.2.8 Tautomerien

Tautomerien sind intramolekulare Umordnungen von Atomen. Dabei werden chemisch zwei verschiedene Moleküle ineinander überführt. Wir wollen dies hier nur exemplarisch am Beispiel der *Keto-Enol-Tautomerie* erklären. Wir betrachten dazu die folgende Abbildung 13

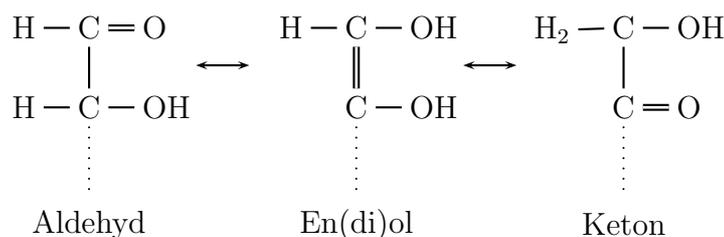


Abbildung 13: Skizze: Keto-Enol-Tautomerie

Im Aldehyd sind aufgrund der Doppelbindung in der Carbonylgruppe und der daraus resultierenden starken Elektronegativität die Elektronen zum Sauerstoffatom der Carbonylgruppe verschoben. Dies führt induktiv zu einer Verlagerung der Elektronen in der C-C Bindung zum Kohlenstoffatom der Carbonylgruppe. Auf der anderen Seite sind die anderen Elektronen im Bindungsorbital zur Hydroxylgruppe aufgrund

derer starken Elektronegativität zum Sauerstoffatom hin verschoben. Dadurch lässt sich das Wasserstoffatom am zweiten Kohlenstoffatom sehr leicht trennen und kann eines der nichtbindenden Orbitale des Sauerstoffatoms der benachbarten Carbonylgruppe protonieren. Diese wandelt sich somit zu einer Hydroxylgruppe und es entsteht zwischen den beiden Kohlenstoffatomen eine Doppelbindung.

Man beachte hierbei, dass die bindenden Orbitale der π -Bindung und die nichtbindenden Orbitale der angrenzenden Sauerstoffatome sich jetzt ebenfalls überlappen, um für die darin enthaltenen Elektronen ein größeres Orbital bereitzustellen. Durch eine Delokalisierung dieser Elektronen kann sich zwischen dem zweiten Kohlenstoffatom und dem Sauerstoffatom der Hydroxylgruppe eine Carbonylgruppe ausbilden. Das frei werdende Wasserstoffatom wird dann unter Aufbruch der Doppelbindung am ersten Kohlenstoffatom angelagert.

Aus ähnlichen Gründen kann sich diese intramolekulare Umlagerung auch auf dem Rückweg abspielen, so dass sich hier ein Gleichgewicht einstellt. Das genaue Gleichgewicht kann nicht pauschal angegeben werden, da es natürlich auch von den speziellen Randbedingungen abhängt. Wie schon erwähnt gibt es auch andere Tautomerien, d.h. intramolekulare Umlagerungen bei anderen Stoffklassen.

0.3 DNS und RNS

In diesem Abschnitt wollen wir uns um die chemische Struktur der *Desoxyribonukleinsäure* oder kurz *DNS* bzw. *Ribonukleinsäure* oder kurz *RNS* (engl. *deoxyribonucleic acid*, *DNA* bzw. *ribonucleic acid*, *RNA*) kümmern. In der DNS wird die eigentliche Erbinformation gespeichert und diese wird durch die RNS zur Verarbeitung weitergegeben. Wie diese Speicherung und Weitergabe im Einzelnen geschieht, werden wir später noch genauer sehen.

0.3.1 Zucker

Ein wesentlicher Bestandteil der DNS sind Zucker. Chemisch gesehen sind Zucker Moleküle mit der Summenformel $C_nH_{2n}O_n$ (weshalb sie oft auch als *Kohlenhydrate* bezeichnet werden). In [Abbildung 14](#) sind die für uns wichtigsten Zucker in der Strukturformel dargestellt. Für uns sind insbesondere Zucker mit 5 oder 6 Kohlenstoffatomen von Interesse. Zucker mit 5 bzw. 6 Kohlenstoffatomen werden auch *Pentosen* bzw. *Hexosen* genannt.

Jeder Zucker enthält eine Carbonylgruppe, so dass Zucker entweder ein Aldehyd oder ein Keton darstellen. Daher werden Zucker entsprechend auch als *Aldose* oder

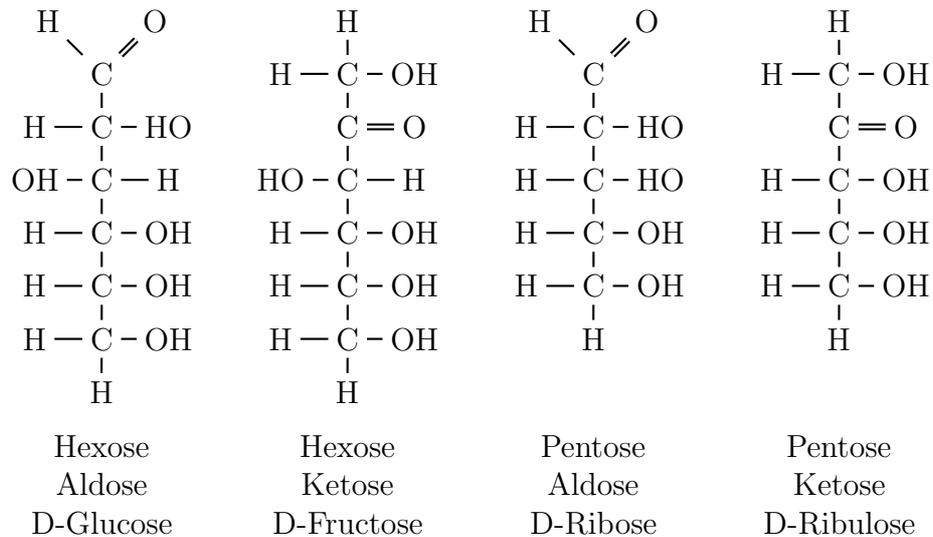


Abbildung 14: Skizze: Zucker (Hexosen und Pentosen sowie Aldosen und Ketosen)

Ketose bezeichnet. Diese Unterscheidung ist jedoch etwas willkürlich, da aufgrund der Keto-Enol-Tautomerie eine Aldose in eine Ketose überführt werden kann. In der Regel pendelt sich ein Gleichgewicht zwischen beiden Formen ein. An dieser Stelle wollen wir noch anmerken, dass es sich bei allen Kohlenstoffatomen (bis auf das erste und das letzte) um asymmetrische Kohlenstoffatome handelt. Somit bilden Zucker Enantiomere aus.

Warum haben jetzt eigentlich Glucose und Fructose unterschiedliche Namen, obwohl diese aufgrund der Keto-Enol-Tautomerie im Gleichgewicht miteinander stehen? In der Natur treten die Zucker kaum als Aldose oder Ketose auf. Die Aldehyd- bzw. Ketogruppe bildet mit einer der Hydroxylgruppen einen Ring aus. In der Regel sind diese 5er oder 6er Ringe. Als 5er Ringe werden diese Zucker als *Furanosen* (aufgrund ihrer Ähnlichkeit zu *Furan*) bezeichnet, als 6er Ringe als *Pyranosen* (aufgrund ihrer Ähnlichkeit zu *Pyran*).

Bei den Hexosen (die hauptsächlich in gewöhnlichen Zuckern und Stärke vorkommen) wird der Ringschluss über das erste und vierte bzw. fünfte Kohlenstoffatom gebildet. Bei Pentosen (mit denen wir uns im Folgenden näher beschäftigen wollen) über das erste und vierte Kohlenstoffatom. Dabei reagiert die Carbonylgruppe mit der entsprechenden Hydroxylgruppe zu einem so genannten *Halb-Acetal*, wie in der folgenden Abbildung 15 dargestellt. Aus der Carbonylgruppe entsteht dabei die so genannte glykosidische OH-Gruppe. Die Ausbildung zum Voll-Acetal geschieht über eine weitere Reaktion (Kondensation) dieser Hydroxylgruppe am zentralen Kohlenstoffatom der ehemaligen Carbonylgruppe.

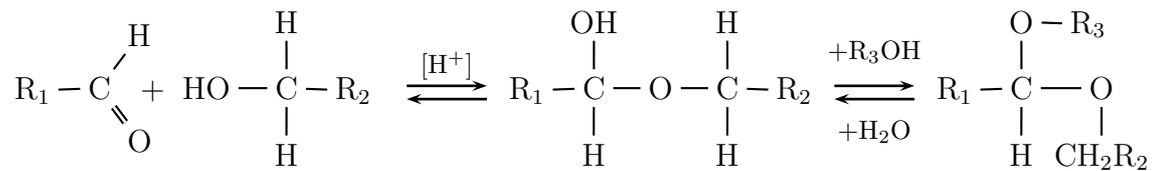


Abbildung 15: Skizze: Halb-Acetal- und Voll-Acetal-Bildung

In der Abbildung 16 sind zwei Furanosen, nämlich *Ribose* und *Desoxyribose* dargestellt. Der einzige Unterschied ist das quasi fehlende Sauerstoffatom am zweiten Kohlenstoffatom, dort ist eine Hydroxylgruppe durch ein Wasserstoffatom ersetzt. Daher stammt auch der Name *Desoxyribose*. Wie man aus dem Namen schon vermuten kann tritt die Desoxyribose in der Desoxyribonukleinsäure (DNS) und die Ribose in der Ribonukleinsäure (RNS) auf. Die Kohlenstoffatome werden dabei zur Unterscheidung von 1 bis 5 durchnummeriert. Aus später verständlich werdenden Gründen, verwenden wir eine gestrichene Nummerierung 1' bis 5' (siehe auch Abbildung 16).

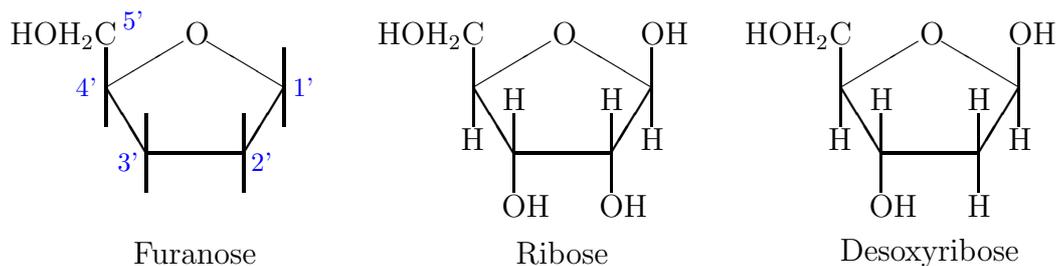


Abbildung 16: Skizze: Ribose und Desoxyribose als Furanosen

0.3.2 Basen

In diesem Abschnitt wollen wir einen weiteren wesentlichen Bestandteil der DNS bzw. RNS vorstellen, die so genannten *Basen*. Hiervon gibt es fünf verschiedene: Adenin, Guanin, Cytosin, Thymin und Uracil. Betrachten wir zuerst die von Purin abgeleiteten Basen. In Abbildung 17 ist links das Purin dargestellt und in der Mitte bzw. rechts *Adenin* bzw. *Guanin*. Die funktionellen Gruppen, die Adenin bzw. Guanin von Purin unterscheiden, sind rot dargestellt. Man beachte auch, dass sich durch die Carbonylgruppe im Guanin auch die Doppelbindungen in den aromatischen Ringen formal ändern. Da es sich hierbei jedoch um alternierende Doppelbindungen und nichtbindende Orbitale handelt, sind die Elektronen sowieso über die aromatischen Ringe delokalisiert.

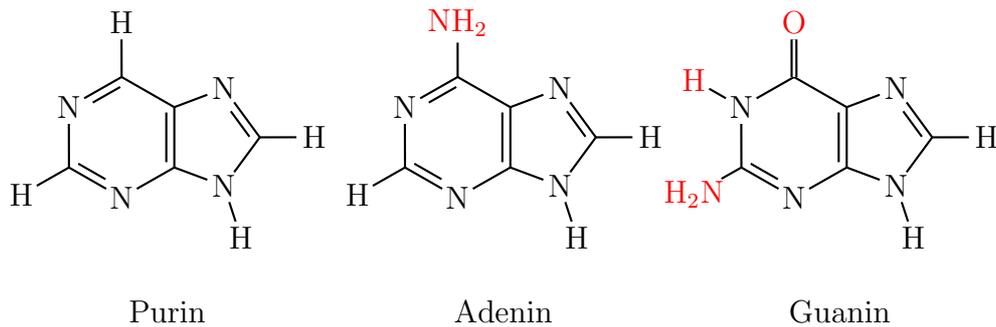


Abbildung 17: Skizze: Purine

Eine weitere Gruppe von Basen erhält man aus Pyrimidin, dessen Strukturformel in der Abbildung 18 links abgebildet ist. Hiervon werden *Cytosin*, *Thymin* und *Uracil* abgeleitet. Auch hier sind die funktionellen Gruppen, die den wesentlichen Unterschied zu Pyrimidin ausmachen, wieder rot bzw. orange dargestellt. Man beachte, dass sich Thymin und Uracil nur in der orange dargestellten Methylgruppe unterscheiden. Hier ist insbesondere zu beachten, dass Thymin nur in der DNS und Uracil nur in der RNS vorkommt. Auch hier beachte man, dass sich durch die Carbonyl-

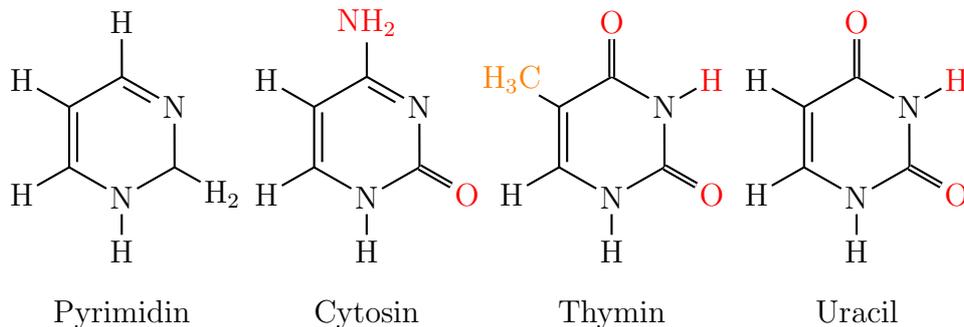


Abbildung 18: Skizze: Pyrimidine

gruppe im Thymin bzw. Uracil auch die Doppelbindungen in den aromatischen Ringen formal ändern. Jedoch bleiben auch hier die Elektronen über die aromatischen Ringe delokalisiert.

Ohne an dieser Stelle im Detail darauf einzugehen, merken wir noch an, dass Guanin über die Keto-Enol-Tautomerie mit einem ähnlichen Stoff in Wechselwirkung steht, ebenso Cytosin über eine so genannte Amino-Imino-Tautomerie. Wir kommen später noch einmal kurz darauf zurück.

Wie im Zucker werden auch die Atome in den aromatischen Ringen durchnummeriert. Da wir im Folgenden auf diese Nummerierung nie zurückgreifen werden, geben

wir sie an dieser Stelle auch nicht an. Um eine Verwechslung mit der Nummerierung in den Zuckern zu vermeiden, wurde die Nummerierung in den Zuckern gestrichen durchgeführt.

0.3.3 Polymerisation

Nun haben wir die wesentlichen Bausteine der DNS bzw. RNS kennen gelernt: die Zucker Desoxyribose bzw. Ribose sowie die Basen Adenin, Guanin, Cytosin und Thymin bzw. Uracil. Zusätzlich spielt noch die Phosphorsäure H_3PO_4 eine Rolle. Je ein Zucker, eine Base und eine Phosphorsäure reagieren zu einem so genannten *Nukleotid*, das sich dann seinerseits mit anderen Nukleotiden zu einem Polymer verbinden kann.

Dabei wird das Rückgrat aus der Phosphorsäure und einem Zucker gebildet, d.h. der Desoxyribose bei DNS und der Ribose bei RNS. Dabei reagiert die Phosphorsäure (die eine mehrfache Säure ist, da sie als Donator bis zu drei Wasserstoffatome abgeben kann) mit den Hydroxylgruppen der Zucker zu einer Esterbindung. Eine Bindung wird über die Hydroxylgruppe am fünften Kohlenstoffatom, die andere am dritten Kohlenstoffatom der (Desoxy-)Ribose gebildet. Somit ergibt sich für das Zucker-Säure-Rückgrat eine Orientierung.

Die Basen werden am ersten Kohlenstoffatom der (Desoxy-)Ribose über eine glykosidische Bindung (zum bereits erwähnten Voll-Acetal) angebunden. Eine Skizze eines Teilstranges der DNS ist in Abbildung 19 dargestellt. Man beachte, dass das Rückgrat für alle DNS-Stränge identisch ist. Die einzige Variabilität besteht in der Anbindung der Basen an die (Desoxy-)Ribose. Eine Kombination aus Zucker und Base (also ohne eine Verbindung mit der Phosphorsäure) wird als *Nukleosid* bezeichnet.

0.3.4 Komplementarität der Basen

Zunächst betrachten wir die Basen noch einmal genauer. Wir haben zwei Purin-Basen, Adenin und Guanin, sowie zwei Pyrimidin-Basen, Cytosin und Thymin (bzw. Uracil in der RNS). Je zwei dieser Basen sind *komplementär* zueinander. Zum einen sind Adenin und Thymin komplementär zueinander und zum anderen sind es Guanin und Cytosin. Die Komplementarität erklärt sich daraus, dass diese Paare untereinander Wasserstoffbrücken ausbilden können, wie dies in Abbildung 20 illustriert ist.

Dabei stellen wir fest, dass Adenin und Thymin zwei und Cytosin und Guanin drei Wasserstoffbrücken bilden. Aus energetischen Gründen werden diese Basen immer

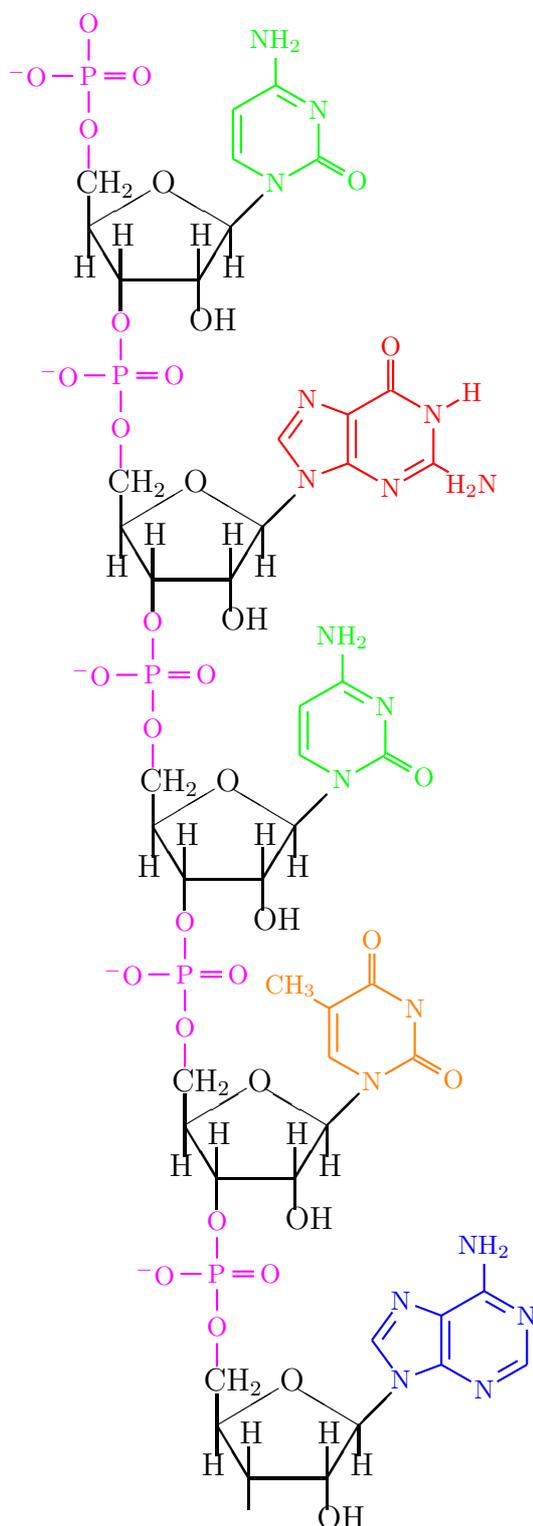


Abbildung 19: Skizze: DNS bzw. RNS als Polymerstrang

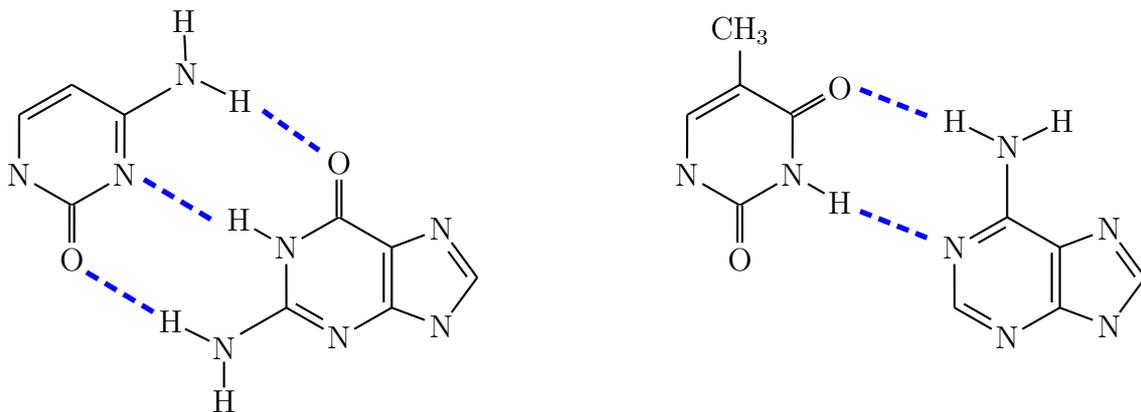


Abbildung 20: Skizze: Wasserstoffbrücken der komplementären Basen

versuchen, diese Wasserstoffbrücken auszubilden. Wir merken, an dass sich auch andere Brückenverbindungen ausbilden können, wie zwischen Thymin und Guanin sowie zwischen Adenin und Cytosin. Diese „falschen“ Wasserstoffbrücken sind aufgrund der Keto-Enol-Tautomerie von Guanin und der Amino-Imino-Tautomerie von Cytosin möglich. Diese sind aus energetischen Gründen zwar eher unwahrscheinlich, können aber dennoch zu Mutationen führen.

Als einfache Merkregel kann man sich merken, dass runde Buchstaben (C und G) bzw. eckige (A und T) zueinander komplementär sind. In der RNS ersetzt Uracil die Base Thymin, so dass man die Regel etwas modifizieren muss. Alles was wie C aussieht ist komplementär (C und **G**) bzw. alles was wie U aussieht (U und **A**).

0.3.5 Doppelhelix

Frühe Untersuchungen haben gezeigt, dass in der DNS einer Zelle die Menge von Adenin und Thymin sowie von Cytosin und Guanin immer gleich groß sind. Daraus kam man auf die Schlussfolgerung, dass diese Basen in der DNS immer in Paaren auftreten. Aus dem vorherigen Abschnitt haben wir mit der Komplementarität aufgrund der Wasserstoffbrücken eine chemische Begründung hierfür gesehen. Daraus wurde die Vermutung abgeleitet, dass die DNS nicht ein Strang ist, sondern aus zwei komplementären Strängen gebildet wird, die einander gegenüber liegen.

Aus sterischen Gründen liegen diese beiden Stränge nicht wie Gleise von Eisenbahnschienen parallel nebeneinander (die Schwellen entsprechen hierbei den Wasserstoffbrücken der Basen), sondern sind gleichförmig miteinander verdreht. Jedes Rückgrat bildet dabei eine Helix (Schraubenlinie) aus. Eine schematische Darstellung ist in [Abbildung 21](#) gegeben.

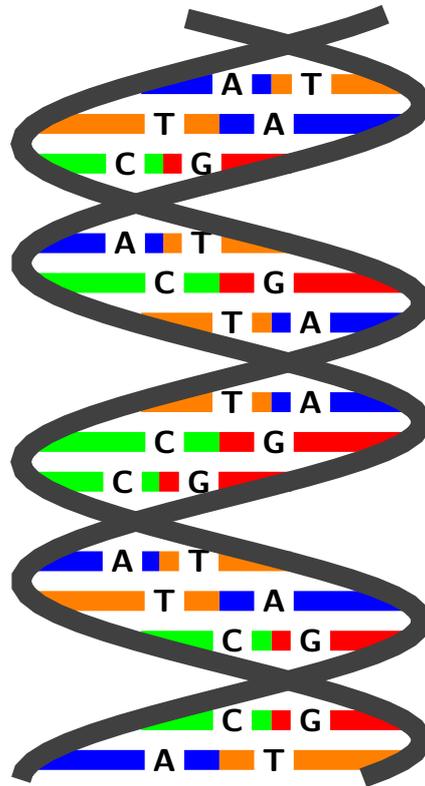


Abbildung 21: Skizze: Doppelhelix der DNS

In einer vollen Drehung sind ungefähr 10 Basenpaare involviert, wobei die Ebenen der Purine bzw. Pyrimidine in etwa orthogonal zur Achse der Doppelhelix liegen. Diese Struktur wurde 1953 von Watson und Crick mit Hilfe der Röntgenkristallographie bestätigt. Es sollte auch noch angemerkt werden, dass unter anderen Randbedingungen auch noch andere Formen von Doppelhelices ausgebildet werden können.

Wie wir schon gesehen haben, besitzt das Rückgrat eines DNS-Strangs eine Orientierung (von 5' nach 3'). Genaue sterische Untersuchungen haben gezeigt, dass die beiden Stränge der DNS innerhalb einer Doppelhelix gegenläufig sind. Läuft also der eine Strang quasi von unten nach oben, so läuft der andere von oben nach unten. Ferner haben die beiden Stränge keinen maximalen Abstand voneinander. Betrachtet man die Doppelhelix der DNS aus etwas „größerem“ Abstand (wie etwa in der schematischen Zeichnung in Abbildung 21), so erkennt man etwas, wie ein kleinere und eine größere Furche auf einer Zylinderoberfläche.

Zum Abschluss noch ein paar Fakten zur menschlichen DNS. Die DNS des Menschen ist nicht eine lange DNS, sondern in 46 unterschiedlich lange Teile zerlegt. Insgesamt sind darin etwa 3 Milliarden Basenpaare gespeichert. Jedes Teil der gesamten DNS ist im Zellkern in einem Chromosom untergebracht. Dazu verdrillt und klumpt

sich die DNS noch weiter und wird dabei von Histonen (spezielle Proteine) unterstützt, die zum Aufwickeln dienen. Würde man die gesamte DNS eines Menschen hintereinander ausrollen, so wäre sie etwa 1 Meter(!) lang.

0.4 Proteine

In diesem Abschnitt wollen wir uns um die wichtigsten Bausteine des Lebens kümmern, die *Proteine*.

0.4.1 Aminosäuren

Zunächst einmal werden wir uns mit den *Aminosäuren* beschäftigen, die den Hauptbestandteil der Proteine darstellen. An einem Kohlenstoffatom (dem so genannten *zentralen Kohlenstoffatom* oder auch *α -ständigen Kohlenstoffatom* ist ein Wasserstoffatom, eine Carboxylgruppe (also eine Säure) und eine Aminogruppe gebunden, woraus sich auch der Name ableitet. Die letzte freie Bindung des zentralen Kohlenstoffatoms ist mit einem weiteren Rest gebunden. Hierfür kommen prinzipiell alle

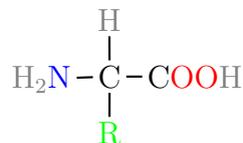


Abbildung 22: Aminosäure

möglichen organischen funktionellen Gruppen in Frage. In der Natur der Proteine kommt jedoch nur eine Auswahl von zwanzig verschiedenen Resten in Betracht. Dabei können die Reste so einfach sein wie ein Wasserstoffatom (Glyzin) oder eine Methylgruppe (Alanin), aber auch recht komplex wie zum Beispiel zwei aromatische Ringe (Tryptophan). In [Abbildung 22](#) ist die Strukturformel einer generischen Aminosäure dargestellt.

In [Abbildung 23](#) sind die Namen der zwanzig in Proteinen auftretenden Aminosäuren und ihre gebräuchlichsten Abkürzungen im so genannten Three-Letter-Code und One-Letter-Code angegeben. Auf die Angabe der genauen chemischen Formeln wollen wir an dieser Stelle verzichten. Hierfür sei auf die einschlägige Literatur verwiesen. In [Abbildung 24](#) sind die grundlegendsten Eigenschaften der einzelnen Aminosäuren schematisch zusammengefasst (nach W.R. Taylor, J. Theor. Biol, 119(2):205–218).

Aminosäure	3LC	1LC
Alanin	Ala	A
Arginin	Arg	R
Asparagin	Asn	N
Asparaginsäure	Asp	D
Cystein	Cys	C
Glutamin	Gln	Q
Glutaminsäure	Glu	E
Glyzin	Gly	G
Histidin	His	H
Isoleuzin	Ile	I
Leuzin	Leu	L
Lysin	Lys	K
Methionin	Met	M
Phenylalanin	Phe	F
Prolin	Pro	P
Serin	Ser	S
Threonin	Thr	T
Tryptophan	Trp	W
Tyrosin	Tyr	Y
Valin	Val	V
Selenocystein	Sec	U
Asparaginsäure oder Asparagin	Asx	B
Glutaminsäure oder Glutamin	Glx	Z
Beliebige Aminosäure	Xaa	X

Abbildung 23: Tabelle: Liste der zwanzig Aminosäuren

Auch hier sind in der Regel (mit Ausnahmen von Glyzin) am zentralen Kohlenstoffatom vier verschiedene Substituenten vorhanden. Somit handelt es sich bei dem zentralen Kohlenstoffatom um ein asymmetrisches Kohlenstoffatom und die Aminosäuren können in zwei enantiomorphen Strukturen auftreten. In der Natur tritt jedoch die L-Form auf, die meistens rechtsdrehend ist! Nur diese kann in der Zelle mit den vorhandenen Enzymen verarbeitet werden.

0.4.2 Peptidbindungen

Auch Aminosäuren besitzen die Möglichkeit mit sich selbst zu langen Ketten zu polymerisieren. Dies wird möglich durch eine so genannte *säureamidartige Bindung* oder auch *Peptidbindung*. Dabei kondensiert die Aminogruppe einer Aminosäure mit der Carboxylgruppe einer anderen Aminosäure (unter Wasserabspaltung) zu einem

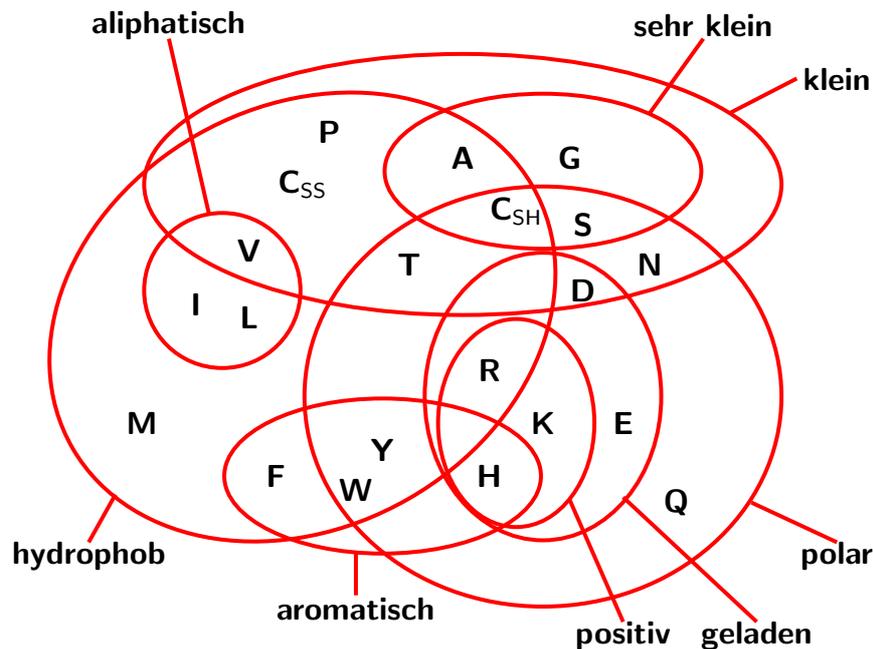


Abbildung 24: Skizze: Elementare Eigenschaften von Aminosäuren

neuen Molekül, einem so genannten *Dipeptid*. Die chemische Reaktionsgleichung ist in Abbildung 25 illustriert.

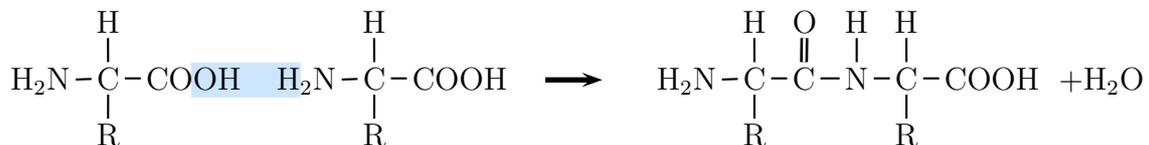


Abbildung 25: Skizze: Säureamidartige oder Peptidbindung

Man beachte, dass das Dipeptid an einem Ende weiterhin eine Aminogruppe und am anderen Ende eine Carboxylgruppe besitzt. Dieser Prozess kann also fortgesetzt werden, so dass sich aus Aminosäuren lange unverzweigte Polymere konstruieren lassen. Solche Polymere aus Aminosäuren nennt man *Polypeptide*. Auch hier bemerken wir wieder, dass ein Polypeptid eine Orientierung besitzt. Wir werden Polypeptide, respektive ihre zugehörigen Aminosäuren immer in der Leserichtung von der freien Aminogruppe zur freien Carboxylgruppe hin orientieren. Ein *Protein* selbst besteht dann aus einem oder mehreren miteinander verbundenen Polypeptiden.

Wir wollen uns nun eine solche Peptidbindung etwas genauer anschauen. Betrachten wir hierzu die Abbildung 26. Von unten links nach oben rechts durchlaufen wir eine Peptidbindung vom zentralen Kohlenstoffatom der ersten Aminosäure über

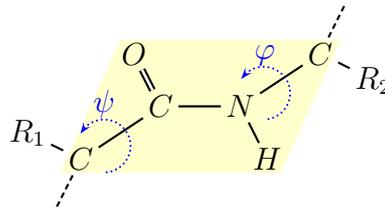


Abbildung 26: Skizze: freie Winkel in der Peptidbindung

das Kohlenstoffatom der ehemaligen Carboxylgruppe über das Stickstoffatom der ehemaligen Aminogruppe der zweiten Aminosäure bis hin zum zentralen Kohlenstoffatom der zweiten Aminosäure.

Auf den ersten Blick könnte man meinen, dass Drehungen um alle drei Bindungen $C-C$, $C-N$ und $N-C$ möglich wären. Eine genaue Betrachtung zeigt jedoch, dass der Winkel um die $C-N$ Bindung nur zwei Werte, nämlich 0° oder 180° , annehmen kann. Dies wird aus der folgenden Abbildung 27 deutlicher, die für die Bindungen $O=C-N$ die beteiligten Elektronen-Orbitale darstellt. Man sieht hier, dass nicht nur die $C=O$ -Doppelbindung ein π -Orbital aufgrund der Doppelbindung ausbildet, sondern dass auch das Stickstoffatom aufgrund seiner fünf freien Außenelektronen zwei davon in einem nichtbindenden p -Orbital unterbringt. Aus energetischen Gründen ist es günstiger, wenn sich das π -Orbital der Doppelbindung und das p -Orbital des Stickstoffatoms überlagern.

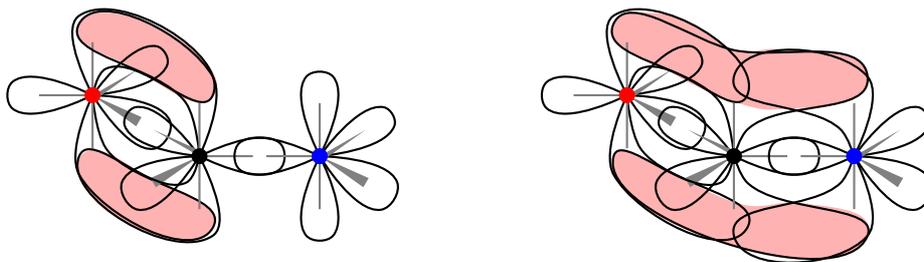


Abbildung 27: Skizze: Elektronenwolken in der Peptidbindung

Somit ist die Bindung zwischen dem Kohlenstoff- und dem Stickstoffatom auf 0° oder 180° festgelegt. In der Regel wird die *trans*-Konformation gegenüber der *cis*-Konformation bevorzugt, da dann die variablen Reste der Aminosäuren ziemlich weit auseinander liegen. Eine Ausnahme stellt nur Prolin dar, da hier die Seitenkette eine weitere Bindung mit dem Rückgrat eingeht.

Prinzipiell unterliegen die beiden anderen Winkel keinen Einschränkungen. Auch hier haben Untersuchungen gezeigt, dass jedoch nicht alle Winkel eingenommen werden. Ein Plot, der alle Paare von den beiden übrigen Winkeln darstellt, ist der so genannte *Ramachandran-Plot*, der schematisch in der Abbildung 28 dargestellt ist. Hier sieht man, dass es gewisse ausgezeichnete Gebiete gibt, die mögliche Winkelkombinationen angeben. Wir kommen auf die Bezeichnungen in diesem Plot später noch einmal zurück.

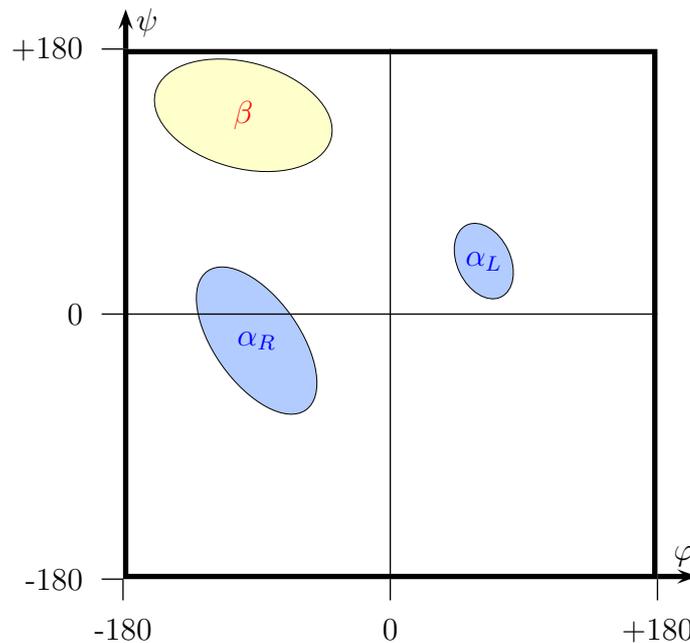


Abbildung 28: Skizze: Ramachandran-Plot (schematische Darstellung)

0.4.3 Proteinstrukturen

Man betrachtet die *Struktur* der Proteine auf vier verschiedenen Ebenen:

0.4.3.1 Primärstruktur

Die *Primärstruktur* (primary structure) eines Proteins ist die Abfolge der beteiligten Aminosäuren des Polypeptids, also seine *Aminosäuresequenz*. Hierbei hält man die Konvention ein, dass man die Aminosäuren von dem Ende mit der freien Amino-Gruppe her aufschreibt. Für uns ist dann ein Protein, respektive seine Primärstruktur nichts anderes als eine Zeichenreihe über einem zwanzig-elementigen Alphabet.

0.4.3.2 Sekundärstruktur

Als Sekundärstruktur (secondary structure) bezeichnet man Regelmäßigkeiten in der lokalen Struktur des Proteins, die sich nur über einige wenige Aminosäuren erstrecken. Die prominentesten Vertreter hierfür sind die spiralförmige α -Helix und der langgestreckte β -Strang (β -strand). Ursache für die Ausbildung dieser Sekundärstrukturmerkmale ist vor allem die Stabilisierung durch Wasserstoffbrückenbindungen. Zu einem großen Teil wird die Sekundärstruktur eines Proteinabschnitts durch seine eigene Primärstruktur bestimmt, d.h. bestimmte Aminosäuresequenzen bevorzugen (oder vermeiden) α -Helices, β -Strands oder Loops.

α -Helices: Wie bei der DNS kann ein Protein oder ein kurzes Stück hiervon eine helixartige (spiralförmige) Gestalt ausbilden. Dabei werden die Helices durch Wasserstoffbrücken innerhalb des Polypeptids stabilisiert, die sich zwischen dem Sauerstoffatom der Carbonylgruppe und dem Wasserstoffatom der Aminogruppe der viertnächsten Aminosäure im Peptidstrang ausbilden. Einen solchen Teil eines Peptids nennt man α -Helix. Dabei entfallen auf eine volle Drehung etwa 3,6 Aminosäuren. Hierbei hat die Helix in der Regel eine Rechtsdrehung, weil bei einer Linksdrehung die sterische Hinderung deutlich größer ist. Die zugehörigen Winkelpaare der Peptidbindung entsprechen im Ramachandran-Plot in Abbildung 28 dem mit α_R markierten Bereich. Einige wenige Helices bilden eine Linksdrehung aus. Die zugehörigen Winkelpaare sind im Ramachandran-Plot mit α_L gekennzeichnet.

Beispielsweise sind die Haare aus Proteinen gebildet, die eine Helix bilden. Auch hier können wie bei der DNS mehrere (sogar mehr als zwei) Helices zusammen verdrillt sein. Ebenso seien Proteine erwähnt, die in Muskeln eine wichtige Rolle spielen.

π - und 3_{10} -Helices: In seltenen Fällen treten Abwandlungen der α -Helix auf, bei denen die Wasserstoffbrücken nicht zwischen den Aminosäuren n und $n + 4$, sondern zwischen den Aminosäuren n und $n + 3$ oder $n + 5$ gebildet werden. In diesen Fällen ist die Helix also etwas mehr oder etwas weniger verdreht. Man nennt diese beiden Formen die 3_{10} - und die π -Helix. Der Name 3_{10} -Helix entstammt dabei der Tatsache, dass die Helix 3 Aminosäuren pro Umdrehung enthält und dass zwischen den beiden Enden einer Wasserstoffbrücke zehn Atome (incl. Wasserstoffatom) liegen. In dieser Nomenklatur (nach Linus Pauling und Robert Corey) würde man die α -Helix mit 3.6_{13} und die π -Helix mit 4.4_{16} bezeichnen.

β -Strands: Eine andere, häufige Struktur sind langgezogene Bereiche von Aminosäuren. Meist lagern sich hier mehrere Stränge (oft von verschiedenen Polypeptidketten, aber durchaus auch nur von einer einzigen) nebeneinander an und

bilden so genannte β -Sheets oder β -Faltblätter aus. Hierbei ist zu beachten, dass sich diese wie Spaghetti nebeneinander lagern. Dies kann entweder parallel oder antiparallel geschehen (Polypeptide haben ja eine Richtung!). Auch hier werden solche Falblätter durch Wasserstoffbrücken zwischen den Amino- und Carbonylgruppen stabilisiert. Diese Struktur heißt Falblatt, da es entlang eines Polypeptid-Strangs immer wieder auf und ab geht, ohne insgesamt die Richtung zu ändern. Bilden sich ganze β -Faltblätter aus, so sehen diese wie ein gefaltetes Blatt aus, wobei die einzelnen Polypeptide quer zu Faltungsrichtung verlaufen. Die zugehörigen Winkelpaare der Peptidbindung entsprechen im Ramachandran-Plot in Abbildung 28 dem mit β markierten Bereich. Beispielsweise tauchen im Seidenfibroin (Baustoff für Seide) fast nur β -Faltblätter auf.

Reverse Turns: Zum Schluss seien noch kurze Sequenzen (von etwa fünf Aminosäuren) erwähnt, die einfach nur die Richtung des Polypeptids umkehren. Diese sind beispielsweise in antiparallelen β -Faltblätter zu finden, um die einzelne β -Strands zu einem β -Sheet anordnen zu können.

0.4.3.3 Supersekundärstruktur

Oft lagern sich zwei oder drei Sekundärstrukturelemente zu sogenannten *Motifs* zusammen.

Hairpins Reverse Turns, die zwischen zwei nebeneinander liegenden antiparallelen β -Strands liegen und deshalb die Form einer Haarnadel nachbilden

Coiled coils bestehen aus zwei verdrillten α -Helices und spielen eine wichtige Rolle in Faser-Proteinen

0.4.3.4 Tertiärstruktur

Die *Tertiärstruktur* (tertiary structure) ist die *Konformation*, also die räumliche Gestalt, eines einzelnen Polypeptids. Sie beschreibt, wie die Elemente der Sekundär- und Supersekundärstruktur sich zu so genannten *Domains* zusammensetzen. Hier ist für jedes Atom (oder Aminosäure) die genaue relative Lage zu allen anderen bekannt. Tertiärstrukturen sind insbesondere deshalb wichtig, da für globuläre Proteine die räumliche Struktur für ihre Wirkung wichtig ist (zumindest in fest umschriebenen Reaktionszentren eines Proteins).

0.4.3.5 Quartärstruktur

Besteht ein Protein aus mehreren Polypeptidketten, wie etwa Hämoglobin, dann spricht man von der *Quartärstruktur* (quaternary structure) eines Proteins. Proteine können aus einem einzelnen Polypeptid oder aus mehreren (gleichen oder unterschiedlichen) Polypeptidketten bestehen.

0.5 Der genetische Informationsfluss

Bislang haben wir die wichtigsten molekularbiologischen Bausteine kennen gelernt. Die DNS, die die eigentliche Erbinformation speichert, und sich in der Regel zu den bekannten Chromosomen zusammenwickelt. Über die Bedeutung der zur DNS strukturell sehr ähnlichen RNS werden wir später noch kommen. Weiterhin haben wir mit den Proteinen die wesentlichen Bausteine des Lebens kennen gelernt. Jetzt wollen wir aufzeigen, wie die in der DNS gespeicherte genetische Information in den Bau von Proteinen umgesetzt werden kann, d.h. wie die Erbinformation weitergegeben (vererbt) wird.

0.5.1 Replikation

Wie wird die genetische Information überhaupt konserviert, oder anders gefragt, wie kann man die Erbinformation kopieren? Dies geschieht durch eine Verdopplung der DNS. An einer bestimmten Stelle wird die DNS entspiralisiert und die beiden Stränge voneinander getrennt, was dem Öffnen eines Reißverschlusses gleicht. Diese Stelle wird auch *Replikationsgabel* genannt. Dann wird mit Hilfe der DNS-Polymerase an beiden Strängen die jeweils komplementäre Base oder genauer das zugehörige Nukleotid mit Hilfe der Wasserstoffbrücken angelagert und die Phosphatsäuren bilden mit den nachfolgenden Zuckern jeweils eine Esterbindung zur Ausbildung des eigentlichen Rückgrates aus. Dies ist schematisch in Abbildung 29 dargestellt, wobei die neu konstruierten DNS-Stücke rot dargestellt sind.

Die Polymerase, die neue DNS-Stränge generiert, hat dabei nur ein Problem. Sie kann einen DNS-Strang immer nur in der Richtung vom 5'-Ende zum 3'-Ende synthetisieren. Nach dem Öffnen liegt nun ein Strang in Richtung der Replikationsgabel (in der ja die Doppelhelix entspiralisiert wird und die DNS aufgetrennt wird) in Richtung vom 3'-Ende zum 5'-Ende vor, der andere jedoch in Richtung vom 5'-Ende zum 3'-Ende, da die Richtung der DNS-Stränge in der Doppelhelix ja antiparallel ist.

Der Strang in Richtung vom 3'-Ende zum 5'-Ende in Richtung auf die Replikationsgabel zu, lässt sich jetzt leicht mit der DNS-Polymerase ergänzen, da dann die

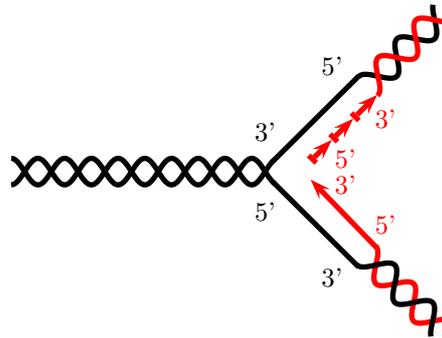


Abbildung 29: Skizze: DNS-Replikation

Synthese selbst in Richtung vom 5'-Ende zum 3'-Ende erfolgt und beim weiteren Öffnen der DNS einfach weiter synthetisiert werden kann.

Für den Strang in Richtung vom 5'-Ende zum 3'-Ende hat man herausgefunden, dass auch hier die Synthese in Richtung vom 5'-Ende zum 3'-Ende erfolgt. Die DNS-Polymerase wartet hier solange, bis ein hinreichend langes Stück frei liegt und ergänzt den DNS Strang dann von der Replikationsgabel weg. Das bedeutet, dass die DNS hier immer in kleinen Stücken synthetisiert wird und nicht im ganzen wie am komplementären Strang. Die dabei generierten kleinen DNS-Teilstränge werden nach ihrem Entdecker *Okazaki-Fragmente* genannt.

0.5.2 Transkription

Damit die in der DNS gespeicherte Erbinformationen genutzt werden kann, muss diese erst einmal abgeschrieben oder kopiert werden. Der Prozess ist dabei im Wesentlichen derselbe wie bei der Replikation. Hierbei wird allerdings nicht die gesamte DNS abgeschrieben, sondern nur ein Teil. Der hierbei abgeschriebene Teil bildet dann nicht eine Doppelhelix mit der DNS aus, sondern löst sich am nichtaktiven Ende wieder, so dass die beiden aufgetrennten Stränge der DNS wieder eine Doppelhelix bilden können.

Hierbei ist zu beachten, dass der abgeschriebene Teil keine DNS, sondern eine RNS ist. Hier wird also Ribose statt Desoxyribose verwendet und als Base Uracil statt Thymin. Die abgeschriebene RNS wird als *Boten-RNS* bzw. *messenger RNA* bezeichnet, da sie als Überbringer der Erbinformation dient.

In prokaryontischen Zellen ist der Vorgang damit abgeschlossen. In eukaryontischen Zellen ist der Vorgang etwas komplizierter, da die Chromosomen im Zellkern beheimatet sind und damit auch die Boten-RNS. Da die Boten-RNS jedoch außerhalb des Zellkerns weiterverarbeitet wird, muss diese erst noch durch die Membran des Zellkerns wandern.

Des Weiteren hat sich in eukaryontischen Zellen noch eine Besonderheit ausgebildet. Die Erbinformation steht nicht kontinuierlich auf der DNS, sondern beinhaltet dazwischen Teilstücke ohne Erbinformation. Diese müssen vor einer Weiterverarbeitung erst noch entfernt werden.

Es hat sich gezeigt, dass dieses Entfernen, *Spleißen* (engl. *Splicing*) genannt, noch im Zellkern geschieht. Dabei werden die Stücke, die keine Erbinformation tragen und *Introns* genannt werden, aus der Boten-RNS herausgeschnitten. Die anderen Teile, *Exons* genannt, werden dabei in der selben Reihenfolge wie auf der DNS aneinander gereiht. Die nach dem Spleißen entstandene Boten-RNS wird dann als *reife Boten-RNS* oder als *mature messenger RNA* bezeichnet.

Für Experimente wird oft aus der mRNA wieder eine Kopie als DNS dargestellt. Diese wird als *cDNS* bzw. *komplementäre DNS* (engl. *cDNA* bzw. *complementary DNA*) bezeichnet. Diese entspricht dann dem Original aus der DNS, wobei die Introns bereits herausgeschnitten sind. Die originalen Gene aus der DNS mit den Introns werden auch als *genetische DNS* (engl. *genetic DNA*) bezeichnet.

Dazu betrachten wir die schematische Darstellung des genetischen Informationsflusses innerhalb einer Zelle (hier einer eukaryontischen) in Abbildung 30.

0.5.3 Translation

Während der *Translation* (*Proteinbiosynthese*) wird die in der DNS gespeicherte und in der reifen Boten-RNS zwischengespeicherte komplementäre Erbinformation in Proteine übersetzt. Dies geschieht innerhalb der Ribosomen, die sich wie zwei Semmelhälften auf den Anfang der Boten-RNS setzen und den RNS-Strang in ein Protein übersetzen. Ribosomen selbst sind aus Proteinen und RNS, so genannter *ribosomaler RNS* oder kurz *rRNS* (engl. *ribosomal RNA*, *rRNA*), zusammengesetzt.

Wir erinnern uns, dass die RNS bzw. DNS im Wesentlichen durch die vier Basen Adenin, Guanin, Cytosin und Uracil bzw. Thymin die Information trägt. Nun ist also die in der RNS gespeicherte Information über einem vierelementigen Alphabet codiert, wobei ein Protein ein Polymer ist, das aus zwanzig verschiedenen Aminosäuren gebildet wird. Wir müssen also noch die Codierung eines zwanzig-elementigen durch ein vier-elementiges Alphabet finden.

Offensichtlich lassen sich nicht alle Aminosäuren durch je zwei Basen codieren. Es muss also Aminosäuren geben, die durch mindestens drei Basen codiert werden. Es hat sich herausgestellt, dass der Code immer dieselbe Länge hat und somit jeweils drei Basen, ein so genanntes *Basen-Triplett* oder *Codon*, jeweils eine Aminosäure codiert.

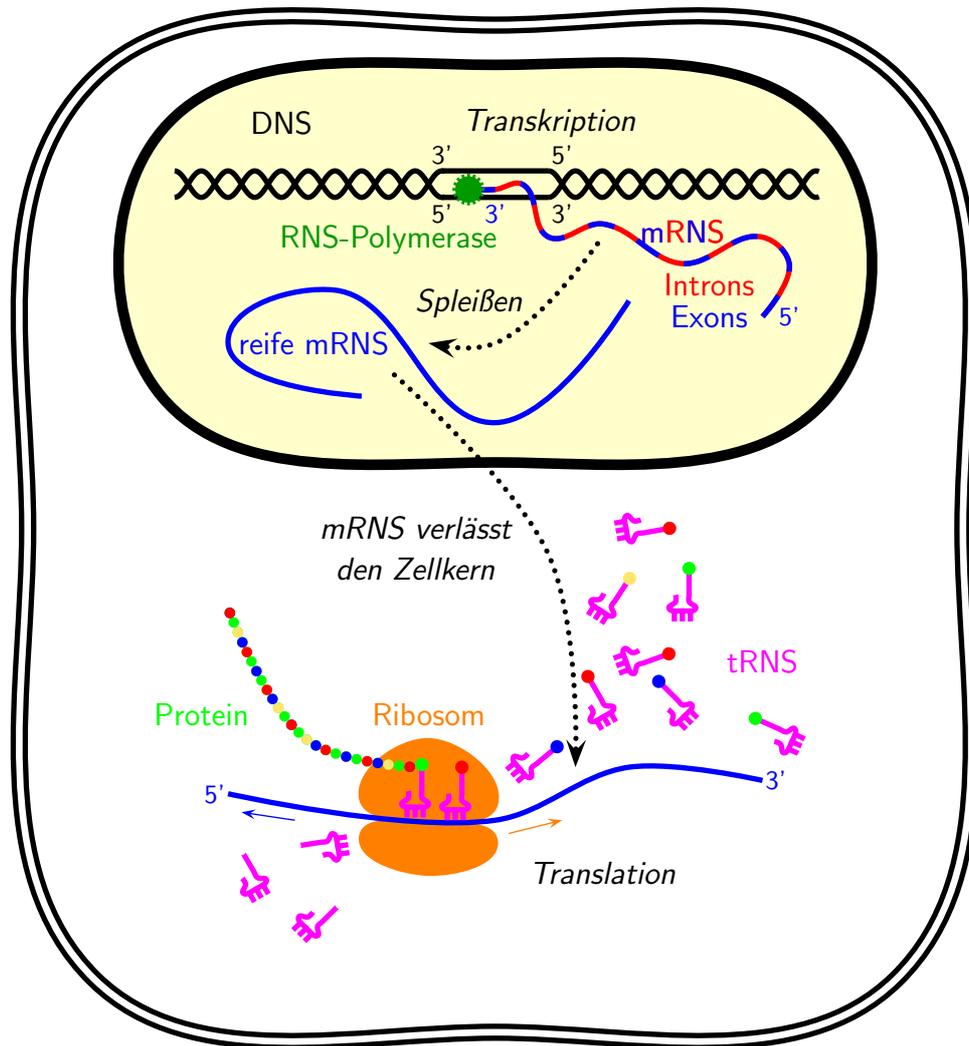


Abbildung 30: Skizze: Transkription und Translation in einer Zelle

Da es 64 verschiedene Triplets aber nur zwanzig Aminosäuren gibt, werden einige Aminosäuren durch mehrere Triplets codiert. In der folgenden Abbildung 31 ist der Code für die Umwandlung von Triplets in Aminosäuren angegeben. In Abbildung 31 steht das erste Zeichen des Triplets links, das zweite oben und das dritte in der rechten Spalte. AGU bzw. AGC codiert also Serin.

Zum genetischen Code ist noch folgendes zu sagen. Es gibt auch spezielle Stopp-Codons, die den Ribosomen mitteilen, dass die Übersetzung der RNS in ein Protein zu Ende ist: Diese sind UAG, UAA und UGA. Ebenfalls gibt es auch ein so genanntes Start-Codon, das jedoch nicht eindeutig ist. AUG übernimmt sowohl die Rolle der Codierung von Methionin als auch dem Anzeigen an das Ribosom, dass hier mit der

	U	C	A	G	
U	Phe	Ser	Tyr	Cys	U
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	STOP	STOP	A
	Leu	Ser	STOP	Trp	G
C	Leu	Pro	His	Arg	U
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	U
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	U
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

Abbildung 31: Tabelle: Der genetische Code

Übersetzung begonnen werden kann. Ebenso ist dieser Code universell, d.h. fast alle Lebewesen benutzen diesen Code. Bislang sind nur wenige Ausnahmen bekannt, die einen anderen Code verwenden, der aber diesem weitestgehend ähnlich ist.

Auch sollte erwähnt werden, dass die Redundanz des Codes (64 Tripletts für 20 Aminosäuren) zur Fehlerkorrektur ausgenutzt wird. Beispielsweise ist die dritte Base für die Decodierung einiger Aminosäuren völlig irrelevant, wie für Alanin, Glycin, Valin und andere. Bei anderen Mutationen werden in der Regel Aminosäuren durch weitestgehend ähnliche (in Bezug auf Größe, Hydrophilie, Ladung oder ähnliches) ersetzt. Auch werden häufig auftretende Aminosäuren durch mehrere Tripletts, selten auftretende nur durch eines codiert.

Ebenfalls sollte man hierbei noch darauf hinweisen, dass es für jeden RNS-Strang eigentlich drei verschiedene Leseraster gibt. Dem RNS-Strang $s_1 \cdots s_n$ an sich sieht man nicht an, ob das codierte Gen in $(s_1 s_2 s_3)(s_4 s_5 s_6) \cdots$, $(s_2 s_3 s_4)(s_5 s_6 s_7) \cdots$ oder $(s_3 s_4 s_5)(s_6 s_7 s_8) \cdots$ codiert ist. Das Start-Codon kann dabei jedoch Hilfe leisten.

Zum Schluss bleibt nur noch die Übersetzung im Ribosom zu beschreiben. Dabei hilft die so genannte *Transfer-RNS* oder *tRNS*. Die Transfer-RNS besteht im Wesentlichen aus RNS und einer Bindungsstelle für eine Aminosäure. Dabei befinden sich drei Basen an exponierter Stelle, die den drei komplementären Basen der zugehörigen Aminosäure entsprechen, die an der Bindungsstelle angebunden ist.

Im Ribosom werden dann jeweils die komplementären tRNS zum betrachteten Triplet der mRNS mittels der Wasserstoffbrücken angebunden. Dabei wird anschließend die neue Aminosäure mittels der säureamidartigen Bindung an den bereits synthetisierten Polypeptid-Strang angebunden. Die tRNS, die den bereits synthetisierten Polypeptid-Strang festhielt, wird dann freigegeben, um in der Zelle wieder mit einer neuen, zum zugehörigen Triplet gehörigen Aminosäure, aufgeladen zu werden.

0.5.4 Das zentrale Dogma

Aus der bisherigen Beschreibung lässt sich die folgende Skizze für den genetischen Informationsfluss ableiten. Die genetische Information wird in der DNS gespeichert und mit Hilfe der Replikation vervielfacht. Mit Hilfe der Transkription wird die genetische Information aus der DNS ausgelesen und in die RNS umgeschrieben. Aus der RNS kann dann mit Hilfe der Translation die genetische Information in die eigentlichen Bausteine des Lebens, die Proteine, übersetzt werden. Damit ist der genetische Informationsfluss eindeutig von der DNS über die RNS zu den Proteinen gekennzeichnet. Dies wird als das zentrale Dogma der Molekularbiologie bezeichnet.

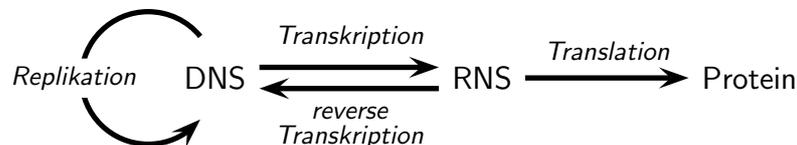


Abbildung 32: Skizze: Das zentrale Dogma der Molekularbiologie

In der Biologie gibt es kaum eine Regel ohne Ausnahmen. Trotz des zentralen Dogmas ist auch ein Informationsfluss in die umgekehrte Richtung möglich. Auch aus der RNS kann modifizierte genetische Erbinformation wieder zurück in die DNS eingebaut werden. Das zentrale Dogma ist in Abbildung 32 noch einmal schematisch dargestellt.

0.5.5 Promotoren

Für den letzten Abschnitt müssen wir uns nur noch überlegen, wie man die eigentliche Erbinformation, die Gene, auf der DNS überhaupt findet. Dazu dienen so genannte *Promotoren*. Dies sind mehrere kurze Sequenzen vor dem Beginn des

eigentlichen Gens, die die RNS-Polymerase überhaupt dazu veranlassen unter bestimmten Bedingungen die spiralisierte DNS an dieser Stelle aufzuwickeln und die beiden DNS-Stränge zu trennen, so dass das Gen selbst transkribiert werden kann. In Bakterien ist dies sehr einfach, da es dort im Wesentlichen nur einen Promotor gibt, der in Abbildung 33 schematisch dargestellt ist. In höheren Lebewesen und insbesondere in eukaryontischen Zellen sind solche Promotoren weitaus komplexer und es gibt eine ganze Reihe hiervon.

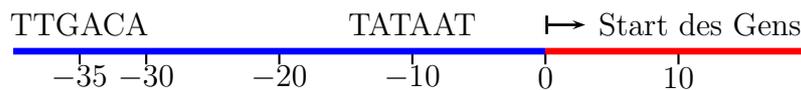


Abbildung 33: Skizze: Promotoren in Bakterien

0.6 Biotechnologie

In diesem Abschnitt wollen wir einige der wichtigsten biotechnologischen Methoden vorstellen, die für uns im Folgenden für eine informatische Modellbildung wichtig sein werden.

0.6.1 Hybridisierung

Mit *Hybridisierung* wird die Aneinanderlagerung zweier DNS-Stränge bezeichnet, die aufgrund der Komplementarität ihrer Basen über Wasserstoffbrücken gebildet wird. Dies haben wir prinzipiell schon bei der Replikation und Transkription kennen gelernt.

Eine Hybridisierung kann dazu ausgenutzt werden, um festzustellen, ob sich eine bestimmte, in der Regel recht kurze Teilsequenz innerhalb eines DNS-Strangs befindet. Dazu wird eine kurze Teilsequenz synthetisiert und an einem Ende markiert, z.B. mit einem fluoreszierenden oder radioaktiven Stoff. Werden die kurzen Teilsequenzen mit den durch Klonierung vervielfachten DNS-Strängen zusammengebracht, können die kurzen Sequenzen mit dem DNS-Strang hybridisieren. Nach Entfernung der kurzen synthetisierten Stücke kann dann festgestellt werden, ob die langen DNS-Stränge fluoreszent oder radioaktiv sind. Letzteres ist genau dann der Fall, wenn eine Hybridisierung stattgefunden hat.

0.6.2 Klonierung

Für biologische Experimente wird oft eine Vielzahl von identischen Kopien eines DNS-Stückes benötigt. Solche Vervielfältigungen lassen sich mit Hilfe niederer Organismen erledigen. Dazu wird das zu vervielfältigende DNS-Stück in die DNS des Organismus eingesetzt und dieser vervielfältigt diese wie seine eigene DNS. Je nachdem, ob Plasmide, Bakterien oder Hefe (engl. yeast) verwendet werden, spricht man vom *plasmid (PAC)*, *bacterial (BAC)* oder *yeast artificial chromosomes (YAC)*.

Die bei PACs verwendeten Plasmide sind ringförmige DNS-Stränge, die in Bakterien auftreten. Bei jeder Zellteilung wird dabei auch der zu klonierende, neu eingesetzte DNS-Strang vervielfältigt. Hierbei können jedoch nur Stränge bis zu 15.000 Basenpaaren kloniert werden.

Bei BACs werden Phagen (ein Virus) verwendet. Die infizierten Wirtszellen (Bakterien) haben dann das zu klonierende DNS-Stück, das in die Phage eingesetzt wurde, vervielfältigt. Hier sind Vervielfältigungen von bis zu 25.000 Basenpaaren möglich.

Bei YACs wird die gewöhnliche Brauerhefe zur Vervielfältigung ausgenutzt, in die die gewünschten DNS-Teilstücke eingebracht werden. Hierbei sind Vervielfältigungen bis zu 1 Million Basenpaaren möglich.

0.6.3 Polymerasekettenreaktion

Eine andere Art der Vervielfältigung ist mit Hilfe der *Polymerasekettenreaktion* (engl. *polymerase chain reaction*) möglich. Diese hatten wir ja schon bei Replikation von DNS-Doppelhelices kennen gelernt. Wir müssen von dem zu vervielfältigenden Bereich nur die Sequenzen der beiden Endstücke von etwa 10 Basenpaaren kennen. Diese kurzen Stücke werden *Primer* genannt

Zuerst werden die DNS-Stränge der Doppelhelix durch Erhitzen aufgespalten. Dann werden sehr viele komplementäre Sequenzen der Primer zugegeben, so dass sie an die Primer der DNS hybridisieren können. Mit Hilfe der Polymerase werden dann ab den Primern in die bekannte Richtung vom 5'-Ende zum 3'-Ende die Einzelstränge der aufgesplitteten DNS zu einem Doppelstrang vervollständigt. Dies ist in Abbildung 34 schematisch dargestellt. Dabei ist das zu vervielfältigende DNS-Stück grün, die Primer rot und der Rest der DNS grau dargestellt. Die Pfeile geben die Synthetisierungsrichtung der Polymerase an.

Einziges Problem ist, dass bei der ersten Anwendung die Polymerase immer bis zum Ende des Strangs läuft. Es wird also mehr dupliziert als gewünscht. Nun kann man dieses Experiment mehrfach (50 Mal) wiederholen. In jedem Schritt werden dabei



Abbildung 34: Skizze: Polymerasekettenreaktion

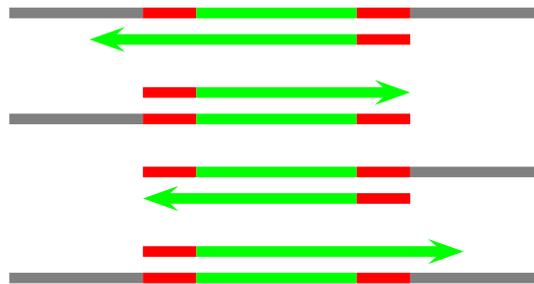


Abbildung 35: Skizze: PCR: nach der 2. Verdopplung

die vorher synthetisierten Doppelstränge verdoppelt. Nach n Phasen besitzt man also 2^n Doppelstränge!

Da nach der ersten Verdopplung bereits jeweils ein unnützes Ende nicht kopiert wurde, überlegt man sich leicht, dass nach der zweiten Verdopplung bereits zwei (von vier) Doppelsträngen nur den gewünschten Bereich verdoppelt haben. Zum Schluss ist jeder zweite DNS-Doppelstrang eine Kopie des gewünschten Bereichs und alle anderen sind nur im gewünschten Bereich doppelsträngig (ansonsten einsträngig, mit zwei Ausnahmen).

Mit dieser Technik lassen sich also sehr schnell und recht einfach eine Vielzahl von Kopien eines gewünschten, durch zwei kurze Primer umschlossenen DNS-Teilstücks herstellen.

0.6.4 Restriktionsenzyme

Restriktionsenzyme sind spezielle *Enzyme* (Proteine, die als Katalysator wirken), die eine DNS-Doppelhelix an bestimmten Stellen aufschneiden können. Durch solche Restriktionsenzyme kann also eine lange DNS geordnet in viele kurze Stücke zerlegt werden. Eines der ersten gefundenen Restriktionsenzyme ist EcoRI, das im Bakterium *Escherichia Coli* auftaucht. Dieses erkennt das Muster GAATTC. In Abbildung 36 ist dieses Muster in der Doppelhelix der DNS noch einmal schematisch mit den Bruchstellen dargestellt. Man beachte hierbei, dass die Sequenz zu sich selbst komplementär ist. Es handelt sich also um ein so genanntes *komplementäres Palindrom*.

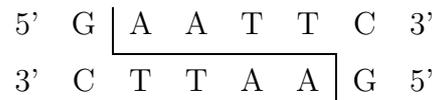


Abbildung 36: Skizze: Restriktionsenzym mit Muster *GAATTC*

0.6.5 Sequenzierung kurzer DNS-Stücke

In diesem Abschnitt wollen wir kurz die gebräuchlichsten Techniken zur Sequenzierung von DNS darstellen. Mit *Sequenzierung* ist das Herausfinden der Abfolge der vier verschiedenen Basen in einem gegebenen DNS-Strang gemeint.

Die Grundidee ist die Folgende: Zuerst wird für das zu sequenzierende Teilstück mit Hilfe der Klonierung eine ausreichende Anzahl identischer Kopien erzeugt. Dann werden die klonierten Teilstücke in vier Gruppen (quasi für jede Base eine) eingeteilt. In jeder Gruppe werden an jeweils einer Base die Teilstücke aufgebrochen. Zu Beginn hat man eines der Enden mit Hilfe eines fluoreszierenden oder radioaktiven Stoffes markiert. Im Folgenden interessieren nur die Bruchstücke, die das markierte Ende besitzen, wobei die anderen Bruchstücke jedoch nicht entfernt werden. Mit Hilfe der so genannten *Elektrophorese* werden die verschieden langen Bruchstücke getrennt.

Die Elektrophorese nutzt aus, dass unter bestimmten Randbedingungen die DNS-Bruchstücke nicht elektrisch neutral, sondern elektrisch geladen sind. Somit kann man die DNS-Bruchstücke mit Hilfe eines elektrischen Feldes wandern lassen. Dazu werden diese innerhalb eines Gels gehalten, dessen Zähflüssigkeit es erlaubt, dass sie sich überhaupt, aber auch nicht zu schnell bewegen. Da die Bruchstücke alle dieselbe Ladung tragen, aber aufgrund ihrer Länge unterschiedlich schwer sind, haben sie innerhalb des angelegten elektrischen Feldes eine unterschiedliche Wanderungsgeschwindigkeit. Die kurzen wandern naturgemäß sehr schnell, während die langen Bruchstücke sich kaum bewegen.

Führt man dieses Experiment gleichzeitig für alle vier Gruppen (also für jede Base) getrennt aus, so erhält man ein Bild der gewanderten Bruchstücke. Dies ist schematisch in Abbildung 37 dargestellt, das man sich als eine Fotografie einer Gruppe radioaktiv markierter Stücke vorstellen kann (auch wenn dies heute mit fluoreszierenden Stoffen durchgeführt wird). Hier ist links die (noch nicht bekannte) Sequenz der einzelnen Bruchstücke angegeben. Mit rot ist speziell das Ergebnis für die Base Adenin hervorgehoben. In den experimentellen Ergebnissen gibt es an sich jedoch keine farblichen Unterscheidungen.

Man kann nun die relativen Positionen einfach feststellen und anhand der Belichtung des Films feststellen in welcher Gruppe sich ein Bruchstück befindet und somit die Base an der entsprechenden Position ablesen. Es ist hierbei zu berücksichtigen, dass

ablesen (siehe auch Abbildung 38). Diese Methode wird nach ihren Erfinder auch *Maxam-Gilbert-Methode* genannt.

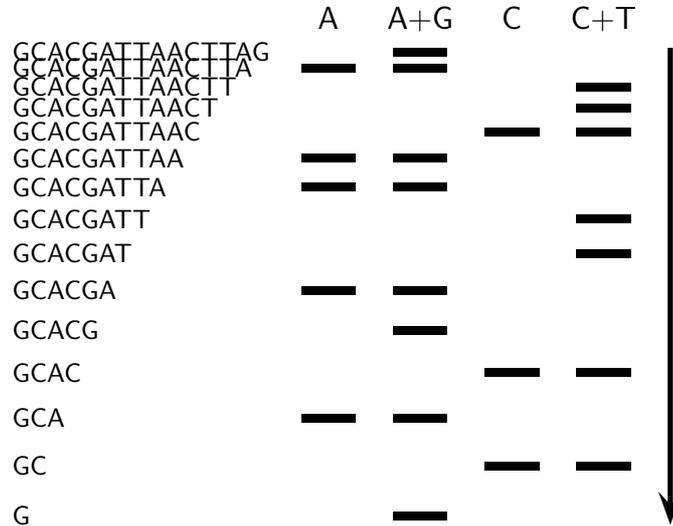


Abbildung 38: Skizze: Sequenzierung nach Maxam-Gilbert

Heutzutage wird in der Regel die Sanger-Methode mit fluoreszierenden Stoffen angewendet, die sich dann gleichzeitig in großen Sequenzierautomaten sehr leicht automatisch anwenden lässt. Dennoch lassen sich auch heute nur DNS-Sequenzen der Länge 500 gut sequenzieren. Mit diesen bekannten Methoden sind auch in Zukunft bei einem entsprechenden technologischen Fortschritt Sequenzierungen von mehr als 1000 Basenpaaren nicht denkbar.

0.6.6 Sequenzierung eines Genoms

Im letzten Abschnitt haben wir gesehen, wie sich kurze DNS-Stücke sequenzieren lassen und dass sich diese Methoden nicht auf beliebig lange DNS-Sequenzen ausdehnen lassen. In diesem letzten Abschnitt wollen wir uns überlegen, wie man ein ganzes Genom sequenzieren kann.

0.6.6.1 Primer Walking

Beim *Primer Walking* ist die Idee, dass man die ersten 500 Basen des Genoms sequenziert. Kennt man diese, so kann man am Ende einen möglichst eindeutigen Primer der Länge ca. zwanzig ablesen und mit der Polymerasekettenreaktion, die Sequenz ab dieser Stelle vervielfältigen. Der folgende Sequenzierungsschritt beginnt

daher am Ende (mit einer kleinen Überlappung) des bereits sequenzierten Anfangsstücks.

Dieses Verfahren lässt sich natürlich beliebig wiederholen. Somit läuft man also mit Primern über die Sequenz und sequenziert die DNS Stück für Stück. In der Anwesenheit von Repeats (Wiederholungen) versagt diese Methode, da man innerhalb eines langen Repeats per Definition keinen eindeutigen Primer mehr finden kann.

0.6.6.2 Nested Sequencing

Auch beim *Nested Sequencing* läuft man Stück für Stück über die Sequenz. Immer wenn man eine Sequenz der Länge 500 sequenziert hat, kann man diese mit Hilfe eines Enzyms (Exonuclease) entfernen und mit der restlichen Sequenz weitermachen.

Beide vorgestellten Verfahren, die die Sequenz Stück für Stück sequenzieren, haben den Nachteil, dass sie inhärent sequentiell und somit bei großen Genomen sehr langsam sind.

0.6.6.3 Sequencing by Hybridization

Beim *Sequenzieren durch Hybridisierung* oder kurz *SBH* (engl. sequencing by hybridization) werden die Sequenzen zuerst vervielfältigt und dann durch Restriktionsenzyme kleingeschnitten. Diese klein geschnittenen Sequenzen werden durch Hybridisierung mit allen Sequenzen der Längen bis etwa 8 verglichen. Somit ist bekannt, welche kurzen Sequenzen in der zu sequenzierenden enthalten sind. Aus diesen kurzen Stücken lässt sich dann mit Informatik-Methoden die Gesamtsequenz wiederherstellen.

Dieses Verfahren ist jedoch sehr aufwendig und hat sich bislang nicht durchgesetzt. Jedoch hat es eine bemerkenswerte Technik, die so genannten *DNA-Microarrays* oder *Gene-Chips* hervorgebracht, die jetzt in ganz anderen Gebieten der Molekularbiologie eingesetzt werden.

Ein DNA-Microarray ist eine kleine Glasplatte (etwa 1cm^2), die in etwa 100 mal 100 Zellen aufgeteilt ist. In jeder Zelle wird ein kleines Oligonukleotid der Länge von bis zu 50 Basen aufgebracht. Mit Hilfe der Hybridisierung können nun parallel 10.000 verschiedene Hybridisierungsexperimente gleichzeitig durchgeführt werden. Die zu untersuchenden Sequenzen sind dabei wieder fluoreszent markiert und können nach dem hochparallelen Hybridisierungsexperiment entsprechend ausgewertet werden.

0.6.6.4 Shotgun Sequencing

Eine weitere Möglichkeit, ein ganzes Genom zu sequenzieren, ist das so genannte *Shotgun-Sequencing*. Hierbei werden lange Sequenzen in viele kurze Stücke aufgebrochen. Dabei werden die Sequenzen in mehrere Klassen aufgeteilt, so dass (in der Regel) eine Bruchstelle in einer Klasse mitten in den Fragmenten der anderen Sequenzen liegt.

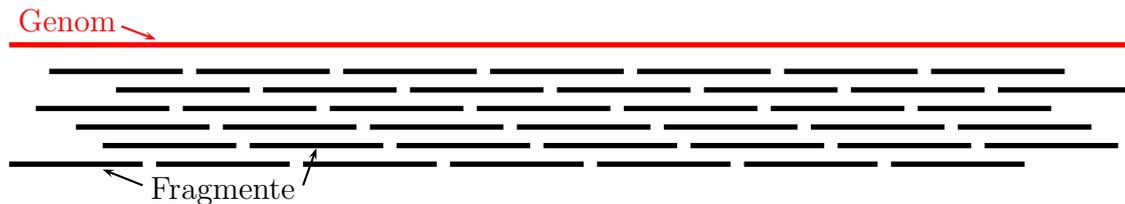


Abbildung 39: Skizze: Shotgun-Sequencing

Die kurzen Sequenzen können jetzt wieder direkt automatisch sequenziert werden. Es bleibt nur das Problem, aus der Kenntnis der Sequenzen wieder die lange DNS-Sequenz zu rekonstruieren. Dabei hilft, dass einzelne Positionen (oder sogar kurze DNS-Stücke) von mehreren verschiedenen Fragmenten, die an unterschiedlichen Positionen beginnen, überdeckt werden. In der Regel sind diese Überdeckungen relativ lang. Somit muss man nur noch die Fragmente wie in einem Puzzle-Spiel so anordnen, dass überlappende Bereiche möglichst gleich sind (man muss ja leider immer noch mit Sequenzierfehlern leben).

Zunächst dachte man, dass diese Methode nur für kürzere DNS-Stränge möglich ist, etwa für 100 000 Basenpaare. Celera Genomics zeigte jedoch mit der Sequenzierung des ganzen Genoms der Fruchtfliege (*Drosophila melanogaster*) und schließlich dem menschlichen Genom, dass diese (bzw. eine geeignet modifizierte) Methode auch für lange DNS-Sequenzen zum Ziel führt.

1.1 Einführendes Beispiel: MSS

In diesem Kapitel befassen wir uns mit einer kurzen Einführung in die Algorithmik. Wir wollen hierzu die Grundlagen zum Entwurf und zur Analyse von Algorithmen besprechen.

Ziel dieses Abschnitts ist es, anhand eines einfachen Problems die verschiedenen Entwurfs- und Analysetechniken für Algorithmen exemplarisch vorzustellen. In den folgenden Abschnitten werden wir dann im einzelnen die benötigten Grundlagen und Techniken genauer behandeln.

1.1.1 Maximal Scoring Subsequence

Zunächst einmal wollen wir das hier betrachtete Problem definieren.

MAXIMAL SCORING SUBSEQUENCE (MSS)

Eingabe: Eine Folge $(a_1, \dots, a_n) \in \mathbb{R}^n$.

Gesucht: Eine (zusammenhängende) Teilfolge (a_i, \dots, a_j) , die $\sigma(i, j)$ maximiert, wobei $\sigma(i, j) := \sum_{\ell=i}^j a_\ell$.

Bemerkung: Mit Teilfolgen sind in diesem Kapitel immer (sofern nicht anders erwähnt) zusammenhängende (d.h. konsekutive) Teilfolgen einer Folge gemeint (also anders als beispielsweise in der Analysis).

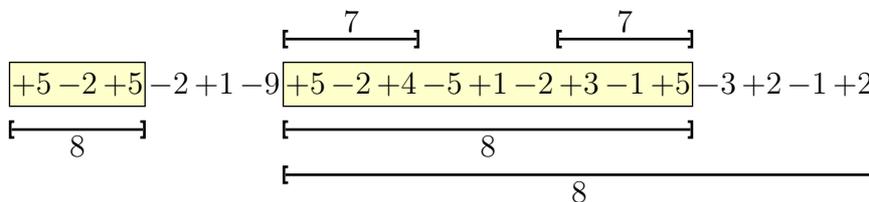


Abbildung 1.1: Beispiel: Maximal Scoring Subsequences

In Abbildung 1.1 ist ein Beispiel angegeben. Wie man dort sieht kann es mehrere (und auch nicht-disjunkte) Lösungen geben.

Bemerkungen:

- Wie man an dem Beispiel sieht, sind mehrere Lösungen möglich.
- Die verschiedenen Lösungen können sich ohne weitere Einschränkung auch überlappen. Ist eine Lösung in der anderen enthalten, so lassen wir im Folgenden als Lösung immer eine Lösung kürzester Länge als echte Lösung zu. Die anderen ergeben sich aus Anhängen von Teilfolgen mit dem Score Null. Darüber hinaus haben solche Lösungen noch eine schöne Eigenschaft, wie wir gleich sehen werden. Dann haben alle Präfixe bzw. Suffixe einer Lösungsfolge einen echt positiven Score.
- Mit der obigen Einschränkung auf kurze Lösungen kann es keine echt überlappenden Lösungen mehr geben. Angenommen, es gäbe echt überlappende Lösungen a' und a'' einer gegebenen Folge a (siehe dazu auch Abbildung 1.2). Die Sequenz gebildet aus der Vereinigung beider Sequenzen (respektive ihrer Indizes) müsste dann einen höheren Score haben, da der Score der Endstücke > 0 ist (sonst würde er in den betrachteten Teilfolgen nicht berücksichtigt werden).

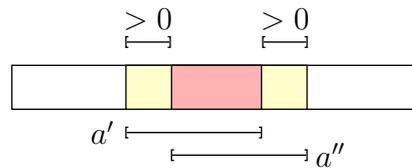


Abbildung 1.2: Skizze: Überlappende optimale Teilfolgen

- Die gleiche Eigenschaft erhält man, wenn man nur längste Lösungen zulässt.

Nur als kurze Motivation wollen wir hier noch anmerken, dass das Problem kein künstliches ist, sondern sogar in der Bioinformatik verschiedene Anwendung besitzt. Beispielsweise kann man diesen Algorithmus zur Bestimmung der *transmembranen Regionen* eines *Transmembranproteins* verwenden. In die Membran eingelagerte Proteine sollten einen ähnlichen Aufbau wie die Membran selbst haben, damit die Gesamtstruktur stabiler ist. Somit sollten transmembrane Regionen hydrophob sein.

Mit einer geeigneten Gewichtung der verschiedenen Aminosäuren können solche hydrophoben Regionen mit Hilfe der Lösung eines Maximal Scoring Subsequence Problems gefunden werden. Für die einzelnen Aminosäuren werden die folgende Werte gemäß der Hydrophobizität der entsprechenden Aminosäure gewählt: Für hydrophobe Aminosäuren wählt man einen positiven Wert aus $[0, 3]$; für hydrophile Aminosäuren einen negativen Wert aus $[-5, 0]$.

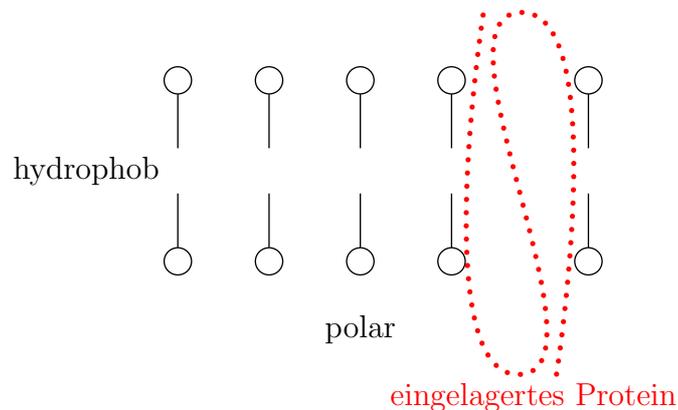


Abbildung 1.3: Beispiel: transmembrane Proteine

Es gibt noch viele weitere Anwendungen, für die wir den interessierten Leser auf das Skript zur *Algorithmen auf Sequenzen* verweisen.

Halten wir noch kurz die eben verwendete Notation formal fest.

Notation 1.1 Für $a, b \in \mathbb{R}$ ist $[a : b] = [a, b] \cap \mathbb{Z} = \{z \in \mathbb{Z} : a \leq z \leq b\}$.

Man beachte, dass wir dabei in der Regel $a, b \in \mathbb{Z}$ verwenden.

1.1.2 Naive Lösung

Im Folgenden wollen wir eine Reihe von Algorithmen zur Lösung des Maximal Scoring Subsequence Problems vorstellen, die jeweils effizienter als die vorher vorgestellte Variante ist.

Beginnen wollen wir mit einem naiven Ansatz. Wir benötigen jedoch vorher noch eine Notation.

Notation 1.2 Sei $a = (a_1, \dots, a_n) \in \mathbb{R}^n$ eine Folge reeller Zahlen, dann bezeichnet $\sigma(i, j) := \sum_{\ell=i}^j a_\ell$ für $i \leq j \in [1 : n]$.

Die naive Methode bestimmt zuerst alle Werte $\sigma(i, j)$ für alle $i \leq j \in [1 : n]$. Eine simple Implementierung ist in Abbildung 1.4 angegeben. Wir merken an, dass wir den Pseudo-Code für die MSS-Algorithmen im Folgenden der Einfachheit halber für ganzzahlige Folgen angeben.

Die Korrektheit des Algorithmus folgt unmittelbar aus der Konstruktion, da wir ja alle Möglichkeiten der Reihe nach ausprobieren. Diese algorithmische Idee nennt

```

MSS_Naive (int[] a, int n)
begin
  maxscore := 0;  ℓ := 1;  r := 0;          /* always is maxscore = σ(ℓ, r) */
  for (i := 1; i ≤ n; i++) do
    for (j := i; j ≤ n; j++) do
      s := 0;                                /* compute s = σ(i, j) */
      for (k := i; k ≤ j; k++) do
        s := s + a[k];
      if (s > maxscore) then
        maxscore := s;  ℓ := i;  r := j;
  end

```

Abbildung 1.4: Algorithmus: naiver Algorithmus für MSS

man auch *vollständige Suche* oder *vollständige Aufzählung* (im Englischen *complete enumeration* oder *exhaustive search*).

Als nächstes stellt sich die Frage, wie gut dieser Algorithmus ist. Was heißt hier überhaupt gut? Zum einen will man natürlich sicher sein, dass der Algorithmus wirklich eine korrekte Lösung liefert. Das heißt, man muss zuerst die *Korrektheit* eines Algorithmus überprüfen. In diesem Fall ist das offensichtlich der Fall.

Zum anderen möchte man natürlich einen möglichst schnellen Algorithmus für das gegebene Problem. Was heißt überhaupt schnell bzw. schneller als was. Dazu muss man also die Laufzeiten von Algorithmen für dasselbe Problem vergleichen. Die einfachste Möglichkeit ist, die gemessenen Laufzeiten für die Implementationen der in Frage kommenden Algorithmen zu vergleichen. Da es sehr viele potenzielle Eingaben gibt, ist das in der Praxis nicht durchführbar, höchstens für eine paar exemplarische Eingaben.

Daher versucht man, die Laufzeit theoretisch abzuschätzen. Eine erste Möglichkeit ist, die Anzahl der Taktzyklen des erzeugten Maschinencodes zu bestimmen. Damit ist man zum einen von der verwendeten Architektur abhängig, zum anderen ist eine solche Bestimmung nicht einfach, da sie in der Regel auch von der Eingabe abhängt.

Ein einfacheres qualitatives Maß ist die Anzahl ausgeführter Befehle des Pseudocodes. Auch dies ist in der Regel sehr aufwendig. Daher ermittelt man meist nur die Anzahl von speziellen ausgeführten Anweisungen, die für den Ablauf des Programms typisch sind. Typisch heißt hier, dass sie die Anzahl aller anderen Operationen majorisieren. Auf jede typische Operation kommt maximal eine konstante Anzahl anderer ausgeführter Operationen.

Dies liefert natürlich keine genaue Vorhersage für die Laufzeit in Sekunden o.ä., aber für einen ersten qualitativen Vergleich von Algorithmen ist das völlig ausreichend, wie wir in diesem Abschnitt noch sehen werden.

Für unser Problem wählen wir als typische Operation die Addition eines Array-Elements. Mit solchen Array-Additionen sind hier und im Folgenden Additionen eines Array-Elements bzw. Additionen einer Variablen, die eine Summe von Array-Elementen enthält, mit einer anderen Zahl (die eventuell ebenfalls ein Array-Element bzw. eine Variable, die eine Summe von Array-Elementen enthält) gemeint.

Damit ergibt sich für unseren Algorithmus der folgende Aufwand $A_{\text{Naiv}}(n)$ für eine Folge mit n Elementen:

$$\begin{aligned}
 A_{\text{Naiv}}(n) &= \sum_{i=1}^n \sum_{j=i}^n \sum_{k=i}^j 1 \\
 &= \sum_{i=1}^n \sum_{j=i}^n (j - i + 1) \\
 &\quad \text{Indexverschiebung in der inneren Summe} \\
 &\quad (\sum_{i=a}^b f(i+c) = \sum_{i=a+c}^{b+c} f(i)) \\
 &= \sum_{i=1}^n \sum_{j=1}^{n-i+1} j \\
 &\quad \text{Gaußsche Summe} \\
 &= \sum_{i=1}^n \binom{(n-i+1)+1}{2} \\
 &\quad \text{Rückwärtige Summation} \\
 &\quad (\sum_{i=a}^b f(c-i) = \sum_{i=c-b}^{c-a} f(i)) \\
 &= \sum_{i=1}^n \binom{i+1}{2} \\
 &\quad \text{Indexverschiebung} \\
 &= \sum_{i=2}^{n+1} \binom{i}{2} \\
 &\quad \text{mit } \sum_{i=k}^n \binom{i}{k} = \binom{n+1}{k+1} \text{ für } k=2 \\
 &= \binom{n+2}{3}
 \end{aligned}$$

$$\begin{aligned}
&= \frac{(n+2)(n+1)n}{3 \cdot 2 \cdot 1} \\
&= \frac{n^3 + 3n^2 + 2n}{6}.
\end{aligned}$$

Lemma 1.3 Sei $(a_1, \dots, a_n) \in \mathbb{R}^n$ eine reelle Folge. Der naive Algorithmus benötigt zur Bestimmung einer maximal scoring subsequence genau $\frac{n^3+3n^2+2n}{6}$ Array-Additionen.

1.1.3 Rekursion

Ein anderes aus Informatik I bekanntes Paradigma zum Entwurf von Algorithmen ist die *Rekursion*. Aufgrund der Assoziativität der Addition gilt folgende Rekursionsgleichung:

$$\sigma(i, j) = \begin{cases} 0 & \text{für } i > j, \\ a_i & \text{für } i = j, \\ \sigma(i, k) + \sigma(k+1, j) & \text{für } i < j \text{ und ein } k \in [i : j-1]. \end{cases}$$

Basierend auf dieser Rekursionsgleichung ergibt sich der in Abbildung 1.5 angegebene Algorithmus.

```

MSS_rec (int[] a, int n)
begin
  maxscore := 0;  ℓ := 1;  r := 0;
  for (i := 1; i ≤ n; i++) do
    for (j := i; j ≤ n; j++) do
      s := σ(i, j);
      if (s > maxscore) then
        maxscore := s;  ℓ := i;  r := j;
      /* compute s = σ(i, j) */
    end
  end
end

int σ(int i, j);
begin
  if (i > j) then return 0;
  else if (i = j) then return a[i];
  else return σ(i, k) + σ(k+1, j);
  /* for some k ∈ [i : j-1] */
end

```

Abbildung 1.5: Algorithmus: rekursive Berechnung des Scores für MSS

Man überlegt sich leicht, dass für die Berechnung Summe $\sigma(i, j)$ mittels Rekursion genau $j - i$ Additionen ausreichend sind. In der Rekursion wird das ursprüngliche Feld der Länge m jeweils in zwei Teile aufgeteilt und die Ergebnisse addiert. Die Aufteilung endet, wenn das Feld einelementig ist. Eine wiederholte Aufteilung eines Feldes in lauter einelementige Felder benötigt genau $m - 1$ Aufteilungen. Also werden in der Rekursion auch genau $m - 1$ Additionen ausgeführt. Damit ergibt sich für diesen Algorithmus der folgende Aufwand $A_{\text{rec}}(n)$ für eine Folge mit n Elementen:

$$\begin{aligned}
 A_{\text{rec}}(n) &= \sum_{i=1}^n \sum_{j=i}^n (j - i) \\
 &= \sum_{i=1}^n \sum_{j=i}^n (j - i + 1) - \sum_{i=1}^n \sum_{j=i}^n 1 \\
 &= A_{\text{Naiv}}(n) - \sum_{i=1}^n (n - i + 1) \\
 &\quad \text{Rückwärtige Summation} \\
 &= A_{\text{Naiv}}(n) - \sum_{i=1}^n i \\
 &\quad \text{Einsetzen von } A_{\text{Naiv}}(n) \text{ und Gaußsche Summe} \\
 &= \frac{n^3 + 3n^2 + 2n}{6} - \frac{n(n+1)}{2} \\
 &= \frac{n^3 + 3n^2 + 2n}{6} - \frac{3n^2 + 3n}{6} \\
 &= \frac{n^3 - n}{6}.
 \end{aligned}$$

Lemma 1.4 Sei $(a_1, \dots, a_n) \in \mathbb{R}^n$ eine reelle Folge. Der auf Rekursion basierte Algorithmus benötigt zur Bestimmung einer maximal scoring subsequence genau $\frac{n^3 - n}{6}$ Array-Additionen.

29.04.25

1.1.4 Dynamische Programmierung

Die im vorherigen Unterabschnitt gefundene Rekursionsgleichung können wir auch noch anders verwenden:

$$\sigma(i, j) = \begin{cases} 0 & \text{für } i > j, \\ a_i & \text{für } i = j, \\ \sigma(i, k) + \sigma(k + 1, j) & \text{für } i < j \text{ und ein } k \in [i : j - 1]. \end{cases}$$

Man überlegt sich leicht, dass bestimmte Summen mehrfach berechnet werden. Mit Hilfe der so genannten *dynamischen Programmierung* können wir jedoch effizienter werden. Wir speichern alle berechneten Werte in einer Tabelle und schauen dort die Werte nach. Dabei wird die Tabelle diagonal von der Mitte nach rechts oben aufgefüllt (siehe auch Abbildung 1.6).

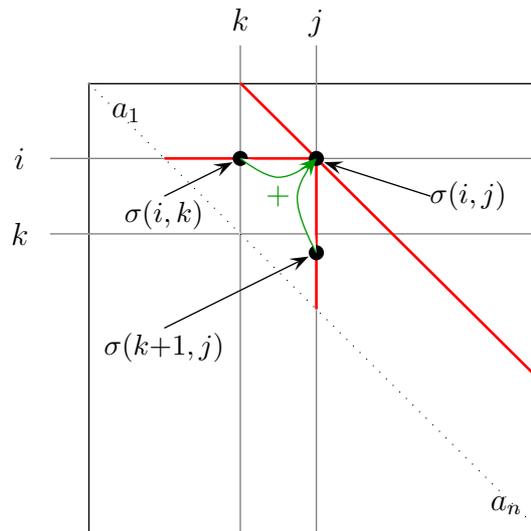


Abbildung 1.6: Skizze: Auffüllen der dynamischen Programmierungstabelle

Wählen wir $k = j - 1$, so können wir die Tabelle auch zeilenweise auffüllen, wie im in Abbildung 1.7 angegebenen Algorithmus. Dabei kann das Feld $S[i, j]$ auch durch eine Variable s ersetzt werden, da immer nur auf das vorhergehende Feldelement in der Zeile zugegriffen wird.

MSS_DP (int[] a , int n)

begin

 maxscore := 0; $\ell := 1$; $r := 0$;

for ($i := 1$; $i \leq n$; $i++$) **do**

for ($j := i$; $j \leq n$; $j++$) **do**

if ($i = j$) **then** $S[i, i] := a[i]$;

else $S[i, j] := S[i, j - 1] + a[j]$;

if ($S[i, j] > \text{maxscore}$) **then**

$\text{maxscore} := S[i, j]$; $\ell := i$; $r := j$;

end

Abbildung 1.7: Algorithmus: Dynamische Programmierung für MSS

Wie ist nun der Aufwand $A_{\text{DP}}(n)$ für die Dynamische Programmierung:

$$\begin{aligned}
 A_{\text{DP}}(n) &= \sum_{i=1}^n \sum_{j=i+1}^n 1 \\
 &= \sum_{i=1}^n (n-i) \\
 &\quad \text{Rückwärtige Summation} \\
 &= \sum_{i=0}^{n-1} i \\
 &= \binom{n}{2} \\
 &= \frac{n^2 - n}{2}.
 \end{aligned}$$

Lemma 1.5 Sei $(a_1, \dots, a_n) \in \mathbb{R}^n$ eine reelle Folge. Der auf dynamischer Programmierung basierende Algorithmus benötigt zur Bestimmung einer maximal scoring subsequence genau $\frac{n^2-n}{2}$ Array-Additionen.

1.1.5 Divide-and-Conquer-Ansatz

Eine andere Lösungsmöglichkeit erhalten wir mit einem *Divide-and-Conquer-Ansatz*, wie in Abbildung 1.8 illustriert. Dabei wird die Folge in zwei etwa gleich lange Folgen aufgeteilt und die Lösung in diesen rekursiv ermittelt.

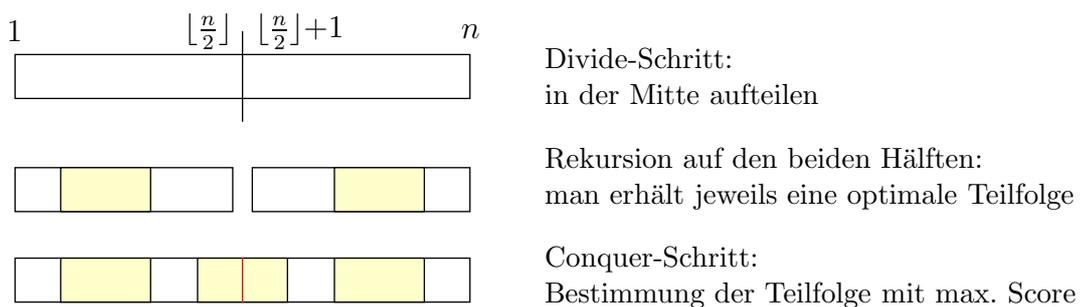
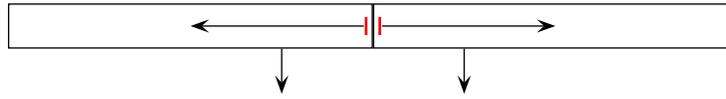


Abbildung 1.8: Skizze: Divide-and-Conquer bei Maximal Scoring Subsequences

Man kann dabei aber auch die optimale Teilfolge in der Mitte zerschneiden. Daher muss man zusätzlich von der Mitte aus testen, wie von dort nach rechts bzw. links eine optimale Teilfolge aussieht (siehe auch Abbildung 1.9).



$$\max \left\{ \sigma(i, \lfloor \frac{n}{2} \rfloor) : i \in [1 : \lfloor \frac{n}{2} \rfloor] \right\} \quad \max \left\{ \sigma(\lfloor \frac{n}{2} \rfloor + 1, j) : j \in [\lfloor \frac{n}{2} \rfloor + 1 : n] \right\}$$

Abbildung 1.9: Skizze: Conquer-Step

Wir halten hier noch schnell die üblichen Definitionen von unteren bzw. oberen Gauß-Klammern und ein paar elementare Eigenschaften fest.

Notation 1.6 Für $x \in \mathbb{R}$ ist

- $\lfloor x \rfloor = \max \{z \in \mathbb{Z} : z \leq x\}$ (untere Gauß-Klammer)
- $\lceil x \rceil = \min \{z \in \mathbb{Z} : z \geq x\}$ (obere Gauß-Klammer)

Beobachtung 1.7 Sei $z \in \mathbb{Z}$, dann gilt $\lfloor z/2 \rfloor + \lceil z/2 \rceil = z$.

Beobachtung 1.8 Sei $x \in \mathbb{R}$, dann gilt $\lfloor x \rfloor = -\lceil -x \rceil$ bzw. $-\lfloor x \rfloor = \lceil -x \rceil$.

Um die optimale Teilfolge zu bestimmen, die die Positionen $\lfloor n/2 \rfloor$ und $\lfloor n/2 \rfloor + 1$ enthält, bestimmen wir zunächst je eine optimale Teilfolge in der linken bzw. rechten Hälfte, die die Position $\lfloor n/2 \rfloor$ bzw. $\lfloor n/2 \rfloor + 1$ enthält. Dazu bestimmen wir jeweils das Optimum der Hälften, d.h

$$\begin{aligned} i' &:= \operatorname{argmax} \{ \sigma(i, \lfloor n/2 \rfloor) : i \in [1 : \lfloor n/2 \rfloor] \}, \\ j' &:= \operatorname{argmax} \{ \sigma(\lfloor n/2 \rfloor + 1, j) : j \in [\lfloor n/2 \rfloor + 1 : n] \}. \end{aligned}$$

Notation 1.9 Sei $f : M \rightarrow N$ eine Abbildung von einer Menge M in eine total geordnete Menge N . Dann ist $\operatorname{argmax} \{f(i) : i \in M\}$ (und analog argmin) definiert als ein Indexwert $j \in M$ mit $f(j) = \max \{f(i) : i \in M\}$.

Wir zeigen jetzt, dass die optimale Teilfolge, die die Positionen $\lfloor n/2 \rfloor$ und $\lfloor n/2 \rfloor + 1$ überdeckt, aus der Konkatination der beiden berechneten optimalen Teilfolgen in den jeweiligen Hälften bestehen muss.

Lemma 1.10 Sei $a = (a_1, \dots, a_n) \in \mathbb{R}^n$ eine reelle Folge. Sei weiter

$$\begin{aligned} i' &:= \operatorname{argmax} \{ \sigma(i, \lfloor n/2 \rfloor) : i \in [1 : \lfloor n/2 \rfloor] \}, \\ j' &:= \operatorname{argmax} \{ \sigma(\lfloor n/2 \rfloor + 1, j) : j \in [\lfloor n/2 \rfloor + 1 : n] \}. \end{aligned}$$

Dann gibt es keine Teilfolge von a mit einem größeren Score als $\sigma(i', j')$, die die Positionen $\lfloor n/2 \rfloor$ und $\lfloor n/2 \rfloor + 1$ enthält.

Beweis: Wir führen den Beweis durch Widerspruch. Sei also $i^* \in [1 : \lfloor n/2 \rfloor]$ und $j^* \in [\lfloor n/2 \rfloor + 1 : n]$ mit $\sigma(i^*, j^*) > \sigma(i', j')$.

Nach Definition von i' und j' gilt:

$$\sigma(i', \lfloor n/2 \rfloor) \geq \sigma(i^*, \lfloor n/2 \rfloor) \quad \text{und} \quad \sigma(\lfloor n/2 \rfloor + 1, j') \geq \sigma(\lfloor n/2 \rfloor + 1, j^*).$$

Also gilt

$$\begin{aligned} \sigma(i', j') &= \sigma(i', \lfloor n/2 \rfloor) + \sigma(\lfloor n/2 \rfloor + 1, j') \\ &\geq \sigma(i^*, \lfloor n/2 \rfloor) + \sigma(\lfloor n/2 \rfloor + 1, j^*) \\ &= \sigma(i^*, j^*) \\ &> \sigma(i', j') \end{aligned}$$

Dies ist offensichtlich ein Widerspruch. ■

Damit ergibt sich folgender Divide-and Conquer-Algorithmus, der in Abbildung 1.10 angegeben ist.

Wie sieht nun die Laufzeit $A_{\text{DC}}(n)$ für diesen Divide-and-Conquer-Algorithmus aus? Zunächst stellen wir fest, dass für eine Folge der Länge 1 keine Additionen auf Array-Elementen ausgeführt werden, also gilt $A_{\text{DC}}(1) = 0$.

Für die linke Hälfte der Folge bestimmen wir rekursiv eine optimale Teilfolge. Hierfür werden nach Definition $A_{\text{DC}}(\lfloor n/2 \rfloor)$ Additionen von Array-Elementen benötigt. Analog werden für die rechte Hälfte $A_{\text{DC}}(\lceil n/2 \rceil)$ Additionen von Array-Elementen benötigt. Für die Bestimmung der optimalen Teilfolge, die die Positionen $\lfloor n/2 \rfloor$ und $\lfloor n/2 \rfloor + 1$ enthält, sind $(\lfloor n/2 \rfloor - 1) + (\lceil n/2 \rceil - 1) + 1 = n - 1$ Additionen von Array-Elementen nötig. Die letzte plus 1 resultiert aus der Summe der beiden Maxima der linken bzw. rechten Hälfte, die das Element an Position $\lfloor n/2 \rfloor$ bzw. $\lfloor n/2 \rfloor + 1$ beinhalten. Somit erhalten wir für die Laufzeit:

$$A_{\text{DC}}(n) = \begin{cases} 0 & \text{falls } n = 1, \\ A_{\text{DC}}(\lfloor n/2 \rfloor) + A_{\text{DC}}(\lceil n/2 \rceil) + (n - 1) & \text{falls } n > 1. \end{cases}$$

```

MSS (int[] a, int n)


---


begin
  | (maxscore,  $\ell$ ,  $r$ ) := MSS_DC( $a$ , 1,  $n$ );
end

(int,int,int) MSS_DC(int[] a, int  $i$ ,  $j$ );
begin
  | if ( $i = j$ ) then
    | if ( $a[i] > 0$ ) then return ( $a[i]$ ,  $i$ ,  $i$ );
    | else return ( $0$ ,  $i$ ,  $i - 1$ );
  | else
    |  $m := \lfloor \frac{i+j-1}{2} \rfloor$ ;
    | ( $s_1$ ,  $i_1$ ,  $j_1$ ) := MSS_DC( $a$ ,  $i$ ,  $m$ );
    | ( $s_2$ ,  $i_2$ ,  $j_2$ ) := MSS_DC( $a$ ,  $m + 1$ ,  $j$ );
    |  $i_3 := m$ ;
    |  $s := a[i_3]$ ;
    |  $simax := s$ ;
    | for ( $k := i_3 - 1$ ;  $k \geq i$ ;  $k--$ ) do
      |  $s := s + a[k]$ ;
      | if ( $s > simax$ ) then
        |  $simax := s$ ;  $i_3 := k$ ;
    |  $j_3 := m + 1$ ;
    |  $s := a[j_3]$ ;
    |  $sjmax := s$ ;
    | for ( $k := j_3 + 1$ ;  $k \leq j$ ;  $k++$ ) do
      |  $s := s + a[k]$ ;
      | if ( $s > sjmax$ ) then
        |  $sjmax := s$ ;  $j_3 := k$ ;
    |  $s_3 := simax + sjmax$ ; /*  $s_3 = \sigma(i_3, j_3)$  */
    | if ( $\max\{s_1, s_2, s_3\} = s_1$ ) then return ( $s_1$ ,  $i_1$ ,  $j_1$ );
    | else if ( $\max\{s_1, s_2, s_3\} = s_2$ ) then return ( $s_2$ ,  $i_2$ ,  $j_2$ );
    | else return ( $s_3$ ,  $i_3$ ,  $j_3$ );
  | end

```

Abbildung 1.10: Algorithmus: Divide-and-Conquer-Algorithmus für MSS

Leider lässt sich diese Laufzeit nun schlecht mit den anderen Laufzeiten vergleichen. Wir benötigen also noch eine geschlossene Form dieser *Rekursionsgleichung*. Da beim Rechnen die Gauß-Klammern einige Probleme bereiten, nehmen wir an, dass

n eine Zweierpotenz ist, also $n = 2^k$. Damit hat die Folge und alle resultierenden Teilfolgen eine gerade Länge, mit Ausnahmen der Folgen der Länge 1. Damit können wir leichter rechnen und das Ergebnis für Längen der Folgen 1 kennen wir ja. Damit ergibt sich

$$A_{\text{DC}}(2^k) = \begin{cases} 0 & \text{falls } k = 0 \quad (n = 1), \\ 2A_{\text{DC}}(2^{k-1}) + (2^k - 1) & \text{falls } k > 0 \quad (n > 1). \end{cases}$$

Wir können nun die Rekursionsgleichung auf der rechten Seite erneut einsetzen und erhalten:

$$\begin{aligned} A_{\text{DC}}(2^k) &= 2A_{\text{DC}}(2^{k-1}) + (2^k - 1) \\ &= 2(2A_{\text{DC}}(2^{k-2}) + (2^{k-1} - 1)) + (2^k - 1) \\ &= 2^2A_{\text{DC}}(2^{k-2}) + 2(2^{k-1} - 1) + (2^k - 1) \end{aligned}$$

Da wir nicht viel schlauer geworden sind, setzen wir die Rekursionsgleichung noch einmal ein:

$$\begin{aligned} &= 2^2(2A_{\text{DC}}(2^{k-3}) + (2^{k-2} - 1)) + 2(2^{k-1} - 1) + (2^k - 1) \\ &= 2^3A_{\text{DC}}(2^{k-3}) + 2^2(2^{k-2} - 1) + 2(2^{k-1} - 1) + (2^k - 1) \end{aligned}$$

Die Summanden am Ende schreiben wir als formale Summe:

$$= 2^3A_{\text{DC}}(2^{k-3}) + \sum_{j=0}^2 2^j(2^{k-j} - 1).$$

Nun können wir eine Vermutung anstellen, wie diese Formel nach i -fachen Einsetzen der Rekursion aussieht:

$$A_{\text{DC}}(2^k) = 2^i A_{\text{DC}}(2^{k-i}) + \sum_{j=0}^{i-1} 2^j (2^{k-j} - 1).$$

Wie können wir nun überprüfen, ob wir richtig geraten haben? Hierfür eignet sich am besten ein Beweis mittels vollständiger Induktion.

Beobachtung 1.11 *Es gilt für $i \in [1 : k]$:*

$$A_{\text{DC}}(2^k) = 2^i A_{\text{DC}}(2^{k-i}) + \sum_{j=0}^{i-1} 2^j (2^{k-j} - 1).$$

Beweis: Induktionsanfang ($i = 1$): Offensichtlich ist

$$A_{\text{DC}}(2^k) = 2A_{\text{DC}}(2^{k-1}) + (2^k - 1) = 2^1 A_{\text{DC}}(2^{k-1}) + \sum_{j=0}^{1-1} 2^j (2^{k-j} - 1).$$

Induktionsschritt ($i \rightarrow i + 1$): Es gilt nach Induktionsvoraussetzung:

$$\begin{aligned} A_{\text{DC}}(2^k) &= 2^i A_{\text{DC}}(2^{k-i}) + \sum_{j=0}^{i-1} 2^j (2^{k-j} - 1) \\ &= 2^i (2A_{\text{DC}}(2^{k-i-1}) + (2^{k-i} - 1)) + \sum_{j=0}^{i-1} 2^j (2^{k-j} - 1) \\ &= 2^{i+1} A_{\text{DC}}(2^{k-(i+1)}) + 2^i (2^{k-i} - 1) + \sum_{j=0}^{i-1} 2^j (2^{k-j} - 1) \\ &= 2^{i+1} A_{\text{DC}}(2^{k-(i+1)}) + \sum_{j=0}^{(i+1)-1} 2^j (2^{k-j} - 1) \end{aligned}$$

Also gilt die Formel auch für $i + 1$ und der Induktionsschluss ist vollzogen. ■

In dieser Formel können wir nun $i = k$ setzen und erhalten:

$$\begin{aligned} A_{\text{DC}}(2^k) &= 2^i A_{\text{DC}}(2^{k-i}) + \sum_{j=0}^{i-1} 2^j (2^{k-j} - 1) \\ &= 2^k A_{\text{DC}}(2^{k-k}) + \sum_{j=0}^{k-1} 2^j (2^{k-j} - 1) \\ &= 2^k A_{\text{DC}}(1) + \sum_{j=0}^{k-1} 2^j (2^{k-j} - 1) \\ &= 2^k \cdot 0 + \sum_{j=0}^{k-1} 2^j (2^{k-j} - 1) \\ &= \sum_{j=0}^{k-1} 2^j (2^{k-j} - 1) \end{aligned}$$

Diese Summe können wir nun leicht lösen:

$$= \sum_{j=0}^{k-1} 2^j \cdot 2^{k-j} - \sum_{j=0}^{k-1} 2^j$$

$$= \sum_{j=0}^{k-1} 2^k - \sum_{j=0}^{k-1} 2^j$$

Mit der geometrischen Reihe $\sum_{j=0}^k x^j = \frac{x^{k+1}-1}{x-1}$ für $x \neq 1$:

$$\begin{aligned} &= k2^k - \frac{2^k - 1}{2 - 1} \\ &= k2^k - 2^k + 1. \end{aligned}$$

Mit $2^k = n$ und damit $k = \log(n)$ erhalten wir

$$A_{\text{DC}}(n) = n \log(n) - n + 1.$$

Diese Lösung gilt natürlich nur solche Werte von n , die Zweierpotenzen sind. Dies ist aber kein Problem, da wir jede Folge am Ende mit 0en auffüllen können und für den Algorithmus nur Folgenlängen betrachten können, die Zweierpotenzen sind. Man überlegt sich leicht, dass dadurch die Menge der Lösungen (mit unserer Vereinbarung) nicht verändert wird.

Dies ist natürlich nicht sehr effizient. Wir werden später noch sehen, dass die Einschränkung auf Zweierpotenzen keine große Einschränkung ist, und diese Lösung im Wesentlichen für alle $n \in \mathbb{N}$ korrekt ist.

Lemma 1.12 *Sei $(a_1, \dots, a_n) \in \mathbb{R}^n$ eine reelle Folge. Der Divide-and-Conquer-Algorithmus benötigt zur Bestimmung einer maximal scoring subsequence genau $n \log(n) - n + 1$ Array-Additionen, sofern n eine Zweierpotenz ist.*

Wir merken hier noch an, dass mit \log immer der Logarithmus zur Basis 2 gemeint ist. Andernfalls werden wir dies immer explizit angeben.

30.04.25

A.1 Lehrbücher zur Vorlesung

- S. Aluru (Ed.): *Handbook of Computational Molecular Biology*; Chapman and Hall/CRC, 2006.
- H.-J. Böckenhauer, D. Bongartz: *Algorithmische Grundlagen der Bioinformatik: Modelle, Methoden und Komplexität*, Teubner, 2003.
- Ph. Compeau, P. Pevzner: *Bioinformatics Algorithms — An Active Learning Approach*, Active Learning Publishers, 3rd Ed., 2018.
- P. Clote, R. Backofen: *Introduction to Computational Biology*, John Wiley and Sons, 2000.
- R.C. Deonier, S. Tavare, M.S. Waterman: *Computational Genome Analysis*, Springer, 2005.
- R. Durbin, S. Eddy, A. Krogh, G. Mitchison: *Biological Sequence Analysis*, Cambridge University Press, 1998.
- D. Gusfield: *Algorithms on Strings, Trees, and Sequences — Computer Science and Computational Biology*, Cambridge University Press, 1997.
- N.C. Jones, P.A. Pevzner: *An Introduction to Bioinformatics Algorithms*, MIT Press, 2004.
- V. Mäkinen, D. Belazzougui, F. Cunial, A.I. Tomescu: *Genome-Scale Algorithm Design: Biological Sequence Analysis in the Era of High-Throughput Sequencing*, Cambridge University Press, 2015.
- J.C. Setubal, J. Meidanis: *Introduction to Computational Molecular Biology*, PWS Publishing Company, 1997.
- W.-K. Sung: *Algorithms in Bioinformatics — A Practical Introduction*, CRC Press, 2009.
- M.S. Waterman: *Introduction to Computational Biology: Maps, Sequences, and Genomes*, Chapman and Hall, 1995.

A.2 Lehrbücher zur Bioinformatik

- I. Eidhammer, I. Jonassen, W.R. Taylor: *Protein Bioinformatics — An Algorithmic Approach to Sequence and Structure Analysis*, Jon Wiley and Sons, 2004.
- W.J. Ewens, G.R. Grant: *Statistical Methods in Bioinformatics*, Springer, 2001.
- A. Isaev: *Introduction to Mathematical Methods in Bioinformatics*, Springer, 2004.
- T. Koski: *Hidden Markov Models for Bioinformatics*, Kluwer Academic Publishers, 2001.
- J.S. Liu: *Bayesian Modeling and Computation in Bioinformatics Research*, in: Current Topics in Computational Molecular Biology, T. Jiang, Y. Xu, M.Q. Zhang (Eds.), MIT Press, 2002.
- D.W. Mount: *Bioinformatics — Sequence and Genome Analysis*, Cold Spring Harbor Laboratory Press, 2001.
- E. Ohlebusch: *Bioinformatics Algorithms — Sequence Analysis, Genome Rearrangements, and Phylogenetic Reconstruction*, Oldenbusch Verlag, 2013.
- P.A. Pevzner: *Computational Molecular Biology — An Algorithmic Approach*, MIT Press, 2000.
- V. Sperschneider: *Bioinformatics — Problem Solving Paradigms*, Springer-Verlag, 2008.

A.3 Lehrbücher zur Algorithmik und Komplexität

- G. Ausiello, P. Crescenzi, G. Gambosi, V. Kann, A. Marchetti-Spaccamela, M. Potasi: *Complexity and Approximation: Combinatorial Optimization Problems and Their Approximability*, Springer, 1999.
- T.H. Cormen, C.E. Leiserson, R.L. Rivest, C. Stein: *Introduction to Algorithms*, MIT Press, 2nd edition, 2001.
- M. Crochemore, W. Rytter: *Jewels of Stringology — Text Algorithms*, World Scientific Press, 2003.
- V. Heun: *Grundlegende Algorithmen*, 2. Auflage, Vieweg, 2003.
- J. Kleinberg, E. Tardos: *Algorithm Design*, Addison Wesley, 2005.

- E.W. Mayr, A. Steger, H.J. Prömel: *Lectures on Proof Verification and Approximation Algorithms*, Lecture Notes in Computer Science, Vol. 1367, Springer, 1998.
- B. Smyth: *Computing Pattern in Strings*, Pearson & Addison Wesley, 2003.
- I. Wegener: *Komplexitätstheorie — Grenzen der Effizienz von Algorithmen*, Springer, 2003.

A.4 Lehrbücher zur Algorithmenanalyse

- R.L. Graham, D.E. Knuth O. Patashnik: *Concrete Mathematics: A Foundation of Computer Science*, Addison-Wesley, 1989.
- D.H. Greene, D.E. Knuth: *Mathematics for the analysis of Algorithms*, Birkhäuser, 1990.
- M. Hofri: *Analysis of Algorithms — Computational methods and Mathematical Tools*, Oxford Press, 1995.
- R. Sedgewick, Ph. Flajolet: *An Introduction to the Analysis of Algorithms*, Addison-Wesley Publishing company, 1996.
- A. Steger: *Diskrete Strukturen I*, Springer, 2001.
- Th. Schickinger, A. Steger: *Diskrete Strukturen II*, Springer, 2001.
- R. Sedgewick, Ph. Flajolet: *An Introduction to the Analysis of Algorithms*, Addison Wesley, 1996.
- H.S. Wilf: *generatingfunctionology*, A K Peters, Ltd., 2005.
www.math.upenn.edu/~wilf/DownldGF.html

Symbole

α -Helix, **27**
 α -ständiges Kohlenstoffatom, **22**
 β -strand, **27**
 π -Bindung, **6**
 π -Orbital, **6**
 σ -Bindung, **6**
 σ -Orbital, **5**
 p -Orbital, **5**
 q -Orbital, **5**
 s -Orbital, **5**
 sp -Hybridorbital, **6**
 sp^2 -Hybridorbital, **6**
 sp^3 -Hybridorbital, **5**

A

Adenin, **16**
Akzeptoratom, **7**
Aldose, **14**
Allel, **2**
Aminosäure, **22**
Aminosäuresequenz, **26**
asymmetrisches Kohlenstoffatom, **12**

B

BAC, **36**
bacterial artificial chromosome, **36**
Basen, **16**
Basen-Triplett, **31**
Benzol, **7**
Bindung
 π -Bindung, **6**
 σ -Bindung, **6**
 ionische, **7**
 kovalente, **5**
Boten-RNS, **30**

C

cDNA, **31**
cDNS, **31**
chiral, **12**
Chromosom, **4**
cis-Isomer, **11**
Codon, **31**
complementary DNA, **31**
complete enumeration, **46**
Crossing-Over-Mutation, **4**
Cytosin, **17**

D

delokalisierte π -Elektronen, **7**
deoxyribonucleic acid, **14**
Desoxyribonukleinsäure, **14**
Desoxyribose, **16**
Dipeptid, **24**
Divide-and-Conquer, **51**
DL-Nomenklatur, **13**
DNA, **14**
 complementary, **31**
 genetic, **31**
DNA-Microarrays, **41**
DNS, **14**
 genetische, **31**
 komplementäre, **31**
Domains, **28**
dominant, **3**
dominantes Gen, **3**
Donatoratom, **7**
Doppelhantel, **5**
dynamische Programmierung, **50**

E

Elektrophorese, **38**
Elterngeneration, **1**
Enantiomer, **12**

Enantiomerie, **11**
 enantiomorph, **12**
 Enzym, **37**
 erste Filialgeneration, **1**
 erste Tochtergeneration, **1**
 exhaustive search, **46**
 Exon, **31**

F

Filialgeneration, **1**
 erste, **1**
 zweite, **1**
 Fischer-Projektion, **12**
 funktionelle Gruppe, **11**
 Furan, **15**
 Furanose, **15**

G

Gauß-Klammer
 obere, **52**
 untere, **52**
 Gen, **2, 4**
 dominant, **3**
 rezessiv, **3**
 Gene-Chips, **41**
 genetic DNA, **31**
 genetische DNS, **31**
 Genom, **4**
 Genotyp, **3**
 Guanin, **16**

H

Halb-Acetal, **15**
 heterozygot, **2**
 Hexose, **14**
 homozygot, **2**
 hydrophil, **10**
 hydrophob, **10**
 hydrophobe Kraft, **10**

I

intermediär, **2**
 Intron, **31**

ionische Bindung, **7**

K

Keto-Enol-Tautomerie, **13**
 Ketose, **15**
 Kohlenhydrate, **14**
 Kohlenstoffatom
 α -ständiges, **22**
 asymmetrisches, **12**
 zentrales, **22**
 komplementäre DNS, **31**
 komplementäres Palindrom, **37**
 Komplementarität, **18**
 Konformation, **28**
 Korrektheit, **46**
 kovalente Bindung, **5**

L

linksdrehend, **13**

M

mature messenger RNA, **31**
 Maxam-Gilbert-Methode, **40**
 Maximal Scoring Subsequence, **43**
 Mendelsche Gesetze, **4**
 messenger RNA, **30**
 mischerbig, **2**
 Motifs, **28**
 mRNA, **30**
 MSS, **43**

N

Nested Sequencing, **41**
 nichtbindendes Orbital, **9**
 Nukleosid, **18**
 Nukleotid, **18**

O

obere Gauß-Klammer, **52**
 Okazaki-Fragmente, **30**
 Orbital, **5**
 π -, **6**
 σ -, **5**
 p , **5**

q-, **5**
s, **5**
sp, **6**
*sp*², **6**
*sp*³-hybridisiert, **5**
 nichtbindendes, **9**

P

PAC, **36**
 Palindrom
 komplementäres, **37**
 Parentalgeneration, **1**
 PCR, **36**
 Pentose, **14**
 Peptidbindung, **23**
 Phänotyp, **3**
 plasmid artificial chromosome, **36**
 polymerase chain reaction, **36**
 Polymerasekettenreaktion, **36**
 Polypeptid, **24**
 Primärstruktur, **26**
 Primer, **36**
 Primer Walking, **40**
 Promotoren, **34**
 Protein, **22, 24, 26**
 Proteinbiosynthese, **31**
 Proteinstruktur, **26**
 Pyran, **15**
 Pyranose, **15**

Q

Quartärstruktur, **29**

R

Ramachandran-Plot, **26**
 rechtsdrehend, **13**
 reife Boten-RNS, **31**
 reinerbig, **2**
 Rekursion, **48**
 Rekursionsgleichung, **54**
 Replikationsgabel, **29**
 rezessiv, **3**
 rezessives Gen, **3**

ribonucleic acid, **14**
 Ribonukleinsäure, **14**
 Ribose, **16**
 ribosomal RNA, **31**
 ribosomaler RNS, **31**
 RNA, **14**
 mature messenger, **31**
 messenger, **30**
 ribosomal, **31**
 transfer, **33**
 RNS, **14**
 Boten-, **30**
 reife Boten, **31**
 ribosomal, **31**
 Transfer-, **33**
 rRNA, **31**
 rRNS, **31**
 RS-Nomenklatur, **13**

S

säureamidartige Bindung, **23**
 Sanger-Methode, **39**
 SBH, **41**
 Sequenzieren durch Hybridisierung,
 41
 Sequenzierung, **38**
 Spleißen, **31**
 Splicing, **31**
 Stereochemie, **11**
 Supersekundärstruktur, **28**

T

Tautomerien, **13**
 Tertiärstruktur, **28**
 Thymin, **17**
 Tochtergeneration, **1**
 erste, **1**
 zweite, **1**
 trans-Isomer, **11**
 transfer RNA, **33**
 Transfer-RNS, **33**
 Translation, **31**
 transmembrane Region, **44**

Transmembranproteins, **44**

tRNA, **33**

tRNS, **33**

U

untere Gauß-Klammer, **52**

Uracil, **17**

V

Van der Waals-Anziehung, **9**

Van der Waals-Kräfte, **9**

vollständige Aufzählung, **46**

vollständige Suche, **46**

W

Wasserstoffbrücken, **8**

Y

YAC, **36**

yeast artificial chromosomes, **36**

Z

zentrales Dogma, **34**

zentrales Kohlenstoffatom, **12, 22**

zweite Filialgeneration, **1**

zweite Tochtergeneration, **1**