

## Introduction

### Single nucleotide polymorphisms (SNPs)

Over the past years single nucleotide polymorphisms (SNPs) have become to be considered ideally suited as DNA markers for establishing genetic linkage and therefore as indicators of genetic diseases.<sup>1,2</sup> The information content of SNPs has been estimated to be only a factor three lower than that of microsatellites. SNPs are defined as biallelic polymorphisms if their frequency is higher than one percent in a general population.<sup>3,4</sup> Otherwise they are termed point mutations. SNPs have been estimated to occur at a frequency of one every thousand nucleotides in the human genome and feature low mutation rates.<sup>5,6</sup> They are divided into noncoding and coding SNPs (nondegenerate, degenerate, nonsynonymous and synonymus).

Information about SNPs could be used in different ways in genetic analysis. Firstly, SNPs might be used as common genetic markers in mapping studies. They can be used for whole-genome scans in pedigree-based linkage analysis of families. A map of about 3,000 SNPs might have the same analytical power for this purpose as a map of 800 microsatellite markers. Currently, the latter are the most frequently used type of DNA marker.<sup>7</sup> Secondly, when the genetics of a disease are studied in individuals in a population, rather than in families, the haplotype distributions and linkage disequilibria can be used to map genes by association. For this purpose, it has been estimated that roughly 100,000-500,000 mapped SNPs will be needed.<sup>8</sup> Many processes and diseases seem to be caused or influenced by complex interactions among multiple genes and environmental factors.<sup>9</sup> These include processes involved in development and aging, and common diseases such as diabetes, cancer, cardiovascular and pulmonary disease, neurological diseases, autoimmune diseases, psychiatric disorders, alcoholism, common birth defects, and susceptibility to infectious diseases, teratogens, and environmental agents.<sup>10</sup> Many of the alleles associated with health problems show low penetrance, meaning that only a few of the individuals carrying them will develop

disease or many haplotypes might influence one of a group of diseases such as diabetes, obesity and cardiovascular disease.<sup>11</sup>

Most of the successes to date in identifying the genes associated with diseases inherited in a Mendelian fashion, and the genetic contribution to common diseases, e.g. BRCA1 and 2 for breast cancer, MODY 1, 2 and 3 for type 2 diabetes, and HNPCC for colon cancer, have been of genes with relatively rare penetrant variant alleles.<sup>12</sup> These genes were well suited to discovery by linkage analysis and positional cloning techniques. Nevertheless, the experimental techniques and strategies useful for finding the low penetrance high frequency alleles involved in disease are usually not the same, and are not as well developed.<sup>13</sup> For example, pedigree analysis of families often does not have sufficient power to identify common, weakly contributing loci.<sup>14</sup> The types of association studies that do have the power to identify such loci efficiently require new approaches and scientific resources to make them as robust and powerful as positional cloning.<sup>15</sup> Association studies using a dense SNP map should allow the identification of disease alleles for complex diseases.<sup>16</sup> Determining the relevance of SNPs for certain phenotypes will require comparative studies (termed case-control studies) of thousands of affected and unaffected individuals.<sup>11</sup> The efficiency of such a strategy was also contested.<sup>17</sup> However, a first example of the identification of a gene involved in a complex disease was shown recently. Mutations in the NOD2 gene were identified by a combination of classic linkage analysis and association studies to contribute to Crohn's disease, which is a complex bowel disease with a prevalence of 1 in 1,000 in western countries.<sup>18,19</sup>

Among the resources needed is a genetic map of much higher density than the microsatellite map. In one project the SNP consortium (TSC), an association of pharmaceutical companies and academic research groups, is creating a dense genome-wide map and database of SNPs that will serve as DNA markers for genotyping experiments. A first draft of a SNP map was published together with the Human Genome Project organisation.<sup>20</sup> The SNPs of the TSC were discovered by shotgun sequencing of a panel of 24 ethnically diverse individuals. To this map the Human Genome Project contributed additional SNPs that were discovered by comparing sequences of overlapping large-insert clones. The combined data were integrated in a first map of the human genome containing 1.42 million SNPs providing an average

density of one SNP every 1.9 kilobases. Allele-frequencies of the SNPs were evaluated in independent populations by pooled re-sequencing. It was estimated that around 60,000 of the discovered SNPs fall into exonic regions of the genome and that approximately 85 % of exons are within 5 kilobases of the closest SNP. Simultaneously, the company Celera announced the establishment of a SNP map of the human genome containing 2.8 million SNPs.<sup>21</sup> Celera integrated SNPs of the public domain with its SNPs discovered by comparison of the sequence of the human genome with other resources using computational methods. The public and the private SNP map will be upgraded regularly.

Beyond large-scale genetic studies of complex diseases, SNPs could find their application in pharmacogenetics for the development of personalised and therefore more effective drugs.<sup>22,23</sup> For example, certain SNPs (or their haplotypes) in genes for proteins like membrane receptors significantly influence the binding affinities to their respective drugs and therefore the drug response could vary dramatically from patient to patient. Or as it was demonstrated in another intriguing example, blood levels and consequently efficiency and side effects of drugs are significantly dependent on drug uptake and metabolism. For example, pharmacogenetics was applied for SNPs in genes encoding transport or hepatic enzymes such as the human drug transporter (MDR1) and drug metabolising (e.g. Cyp3A) enzymes.<sup>24</sup>

A further important application of SNPs could be genetic fingerprinting of domestic animals and plants. SNPs are considered to be suitable DNA markers to establish traceability in the agricultural sector. Traceability might become very crucial as it was realised during the BSE crisis in Europe. It is currently not possible to perform full coverage genetic fingerprinting by microsatellite DNA markers because of the required gel-based analysis that is expensive and cumbersome.

### *Techniques for large-scale SNP genotyping*

By most DNA analysis techniques that are amenable for automation, SNPs are more easily analysed than microsatellites.<sup>25</sup> The major current task is the identification of

significant and disease associated SNPs. For this, efficient screening methods are required that allow the analysis of a progressing number of SNP markers on a large number of individuals. In this chapter the most important gel-based, plate-reader and real-time analyse techniques for large-scale SNP genotyping are described.<sup>26-28</sup> However the choice might remain arbitrarily. Currently all of the technologies described in detail below are considered to be too expensive for large-scale SNP genotyping.<sup>29</sup> Further development and improvements will take place. Prices per SNP analysis of the methods described here are around 0.5 US-Dollar and more, while prices between 0.01 and 0.1 US-Dollar are desired.

### *SNP genotyping by DNA sequencing*

DNA sequencing is one of the most important molecular-biological techniques today. It was used to decode the complete genome sequence of different organisms.<sup>30</sup> The applied whole-genome approaches and expressed sequence tag sequencing have started exerting a significant influence on biology and medicine.

Conventional DNA sequencing is based on the enzymatic chain termination method developed more than twenty years ago by F. Sanger and colleagues.<sup>31</sup> Gel electrophoresis is employed for separating the base-specific terminated DNA fragments after a primer-directed polymerase reaction according to their size. Usually, fragment ladders up to 1,000 nucleobases can be analysed. During the last 15 years, many protocols using fluorescently labelled dideoxynucleotides were developed for efficient sequencing of whole genomes.<sup>32</sup> The development of capillary electrophoresis allows even more efficient sequencing because gel-loading can be done automatically and the time for one gel-run was reduced to approximately two hours.<sup>33</sup>

One of the major problems of DNA sequencing is the occurrence of band compressions during electrophoresis of DNA fragments, which have been observed roughly once every 2 kb in human cDNA.<sup>34</sup> Band compressions were observed when a particular DNA fragment migrates faster than expected, leading to an overlap with an adjoining fragment. This posed significant problems in the interpretation of results. For

example, sequence motifs like palindromic sequences of promotor regions could build stable hairpin structures, thereby leading to abnormal migration in electrophoresis. Particularly in these regions that play an important role in cell regular processes, SNP genotyping is interesting but conventional sequencing is often not feasible. Moreover, conventional DNA sequencing is too cost intensive and inefficient for repetitive large-scale SNP genotyping. Kits for SNP genotyping using (mini-) sequencing of approximately 10 bases are commercially available.<sup>35</sup>

An alternative method to conventional DNA sequencing avoiding the described disadvantages is pyrosequencing (figure 1.1).<sup>36</sup> Pyrosequencing is feasible for sequencing approximately 20 bases. This method is based on coupled enzymatic reactions executed in a single tube and furthermore it is a nonelectrophoretic technique. The detection is based on the release of pyrophosphate during DNA polymerase reaction. A DNA fragment like a single-stranded PCR product is incubated with a primer, a DNA polymerase, an ATP sulfurylase, a firefly luciferase and a nucleotide-degrading enzyme (preferentially an apyrase). Repeated cycles of deoxynucleotide addition are performed. Only in the case of complementarity with the nucleotide it will be incorporated into the growing strand.

The synthesis of the growing strand is accompanied by the release of pyrophosphate in equal amounts to that of the incorporated nucleotide. Thereby, real-time signals are acquired by the enzymatic production of inorganic pyrophosphate. The released pyrophosphate is fully converted to ATP by an ATP sulfurylase. The concentration of ATP is then measured by the amount of light produced during a luciferase reaction using suitable light-sensitive detectors such as charge-coupled device (CCD) cameras or luminometers. Unincorporated nucleotides and the produced ATP are degraded between the cycles by an apyrase. The four employed enzymes are synchronised so that sequencing of approximately 20 bases is possible.

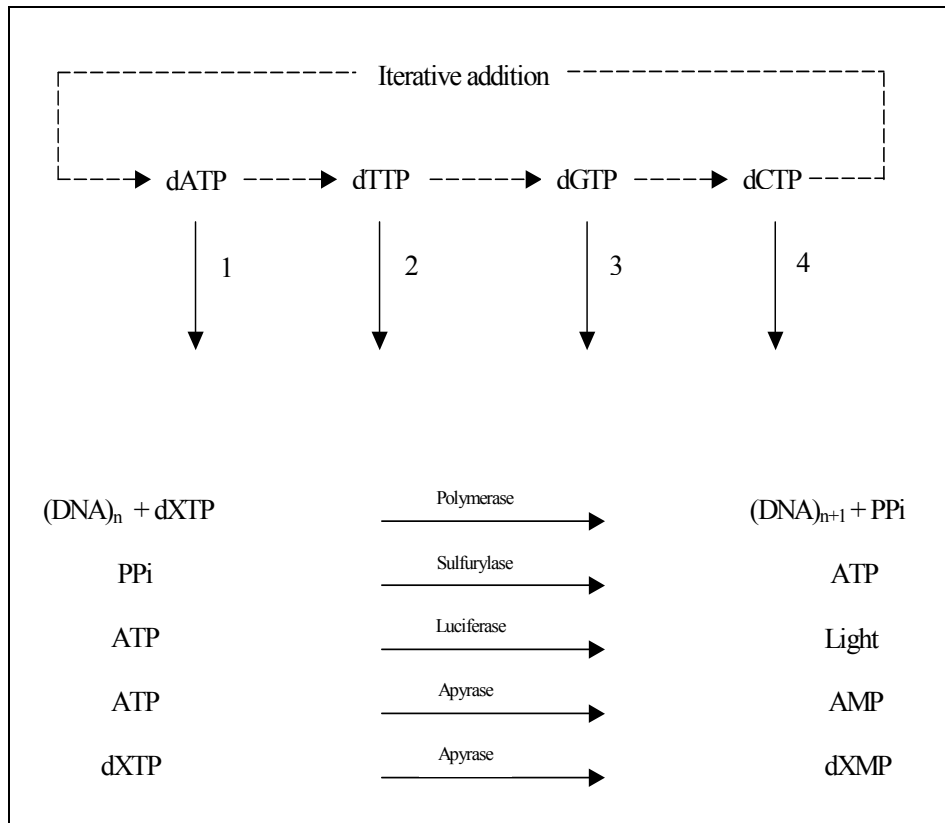


Figure 1.1. The principle of the pyrosequencing procedure. Four nucleotides are added stepwise to the DNA template hybridised to an extension primer. The released pyrophosphate (PPi) is converted to ATP by a sulfurylase and coupled to a luciferase reaction. A nucleotide-degrading enzyme digests the added nucleotides. After the degradation of the first added nucleotide, the next can be added. This procedure can be repeated several times and therefore longer stretches of the template sequence can be deduced. dXTP stands for one of the four nucleotides.

It was claimed by the inventors of pyrosequencing that parallel processing of large numbers of samples would be possible using high-density microtiter plates and microinjector technology. Recently, a device was constructed employing “ink-jet” technology for delivering the nucleotides into the microtiter plates with simultaneous detection of all samples by CCD analysis. The disadvantage of this system seems to be the awkward single-strand preparation from double-stranded PCR products and expensive purification by magnetic beads. This might contribute to the major cost of this technology. A full SNP genotyping system using pyrosequencing is commercially available.<sup>37</sup>

*SNP genotyping by microarrays*

A large number of publications were discerned to an emerging technology for the study of nucleic acid variations called microarrays or DNA chips. Microarrays can be used for (re-) sequencing DNA, for RNA expression studies and for polymorphism detection, particularly for SNP identification and genotyping.<sup>38</sup> The identification of new DNA variations has been executed with oligonucleotides spotted to arrays. The ability to synthesise oligonucleotides at a very high density allowed constructing an (“tiling”) array capable of scanning a target sequence for SNPs. In principle each overlapping 25-mer of the DNA sequence was covered by complementary oligonucleotide probes that differ by carrying an A, T, C or G substituted at each position of the oligonucleotides. Alteration of hybridisation patterns of PCR products that annealed to these probes revealed SNPs. For hybridisation arrays the choice of buffers, hybridisation times and washing conditions is crucial. The optimisation of protocols significantly varies with the sequences that have to be hybridised.<sup>28,29</sup>

Recently, the application of SNPs to large-scale genotyping using DNA arrays has been demonstrated.<sup>39</sup> 100 tiling arrays were used to scan for SNPs in 21,000 sequence tagged sites (STSs) covering roughly 2 Mb of genomic DNA. Currently such arrays implicate the big disadvantage that homozygous variants are detected correctly but heterozygous variants are under-represented. On microarrays only a relatively low signal-background ratio is achieved. Consequently only a small difference in fluorescence upon hybridisation of labelled target oligonucleotides to the matched or mismatched sequence is observed significantly complicating the analysis.<sup>38,39</sup>

Theoretically, rather big regions of the genome could be surveyed for DNA variation. If the actual feature size could be reduced 20-fold to 1 micrometer, it would be possible to investigate 100 Mb on a single array with 4 cm<sup>2</sup> - not to mention the possibility to survey an entire human genome with 30 arrays.<sup>38</sup> Unfortunately, before that several problems have to be solved. First of all, present array technology allows hybridisation of total mammalian RNA but not of genomic DNA with its more than 100-fold higher complexity. Therefore, each target locus requires developing a specific PCR. Secondly, extreme miniaturisation of arrays requires the development of more

sensitive labelling and detection methods. A further problem of hybridisation technologies is the susceptibility to secondary-structure formation.

Interesting variations are electric field (EF) microarrays.<sup>40,41</sup> Instead of chemically binding an oligonucleotide on a surface, electric fields are applied to direct these to a specific address on the support, which is a modified agarose matrix. The hybridisation of templates such as PCR products is also conducted by electric fields resulting in a significant decrease of hybridisation and washing times. The fidelity of this microarray format is increased compared to simple arrays.

In order to improve the specificity of SNP analysis on microarrays, the hybridisation was coupled with an enzymatic step such as primer extension. In an approach termed arrayed primer extension (APEX) PCR products were hybridised to arrayed oligonucleotides.<sup>42-45</sup> In primer extension reactions each primer was extended with a respective fluorescently-labelled dideoxynucleotide by a DNA-polymerase. The polymerase extends the 3'-end of the primer by specifically incorporating nucleotides that are complementary to the DNA template. Extension terminates at the first base in the template where a nucleotide occurs that is complementary to one of the ddNTPs in the reaction mix. The fluorescence emission was therefore specific for the SNP to be analysed on the hybridised PCR target. In a variant of this procedure allele-specific oligonucleotides were used, in which the 3'-ends matched one of the two alleles. Because more than one fluorescent chromophore was incorporated for each matching allele, the detection sensitivity was increased.

Genotyping of large sets of SNPs rests a difficult task. Arrays containing specific detectors for each allele at many loci have been constructed. A more powerful approach seems to be the construction of generic arrays containing "tag sequences".



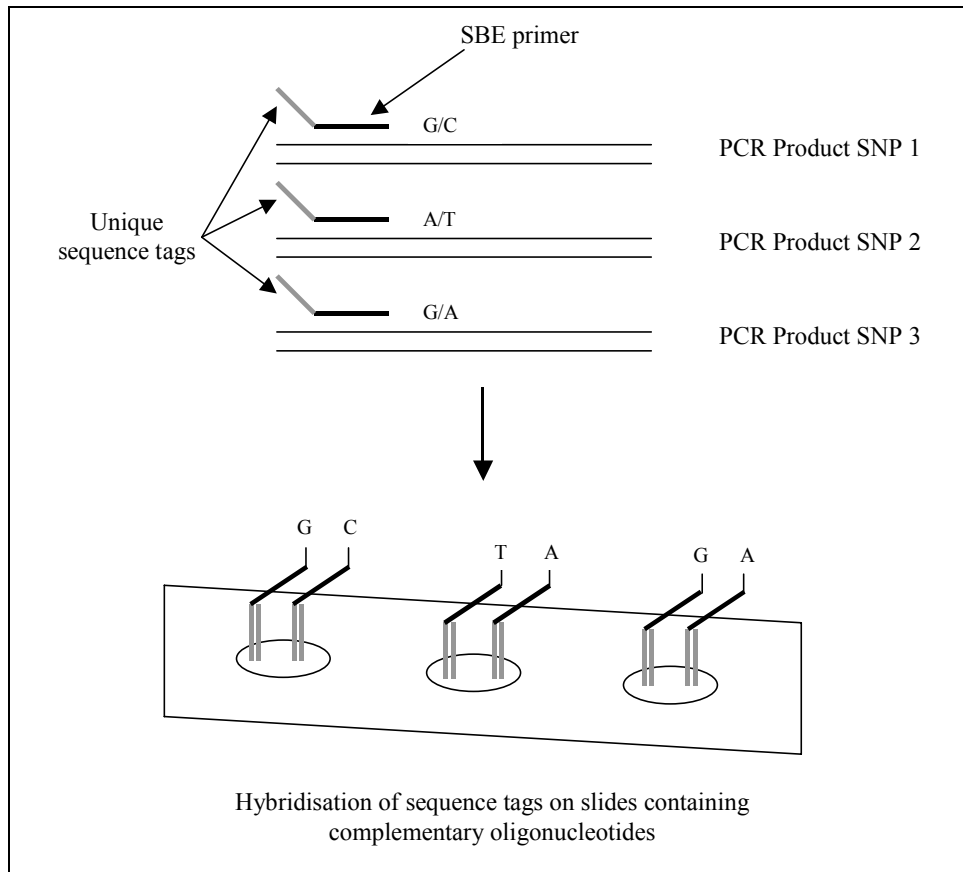


Figure 1.2. The principle of SBE-TAGS genotyping - a promising microarray-based technology for SNP genotyping. In SBE, a primer containing a generic sequence tag and a locus-specific sequence adjacent to the SNP is used for primer extension with fluorescent ddNTPs. Multiple SBE reactions can be performed in solution. Thereby each SBE primer is characteristic having a different unique sequence tag. The multiplexed reaction is analysed after hybridisation to a generic tag array.

If primers for each locus are designed with a unique tag sequence, allele-specific reactions such as primer extension can be done in solution. Then the tag-sequences are hybridised using the same protocol to a generic array so that each assay product anneals to its corresponding sequence. This approach is called single base extension-tag array on glass slides (SBE-TAGS) and can be applied to oligonucleotide and spotted arrays (figure 1.2).<sup>46</sup> This SBE-Tags were developed and are applied at the Whitehead Institute (Cambridge, MA).

Another hybridisation-based technology coupled with an allele-specific step is fluorescence-labelled, coded microspheres.<sup>47,48</sup> The microsphere technology for SNP genotyping is commercially available.<sup>49</sup>

## Introduction

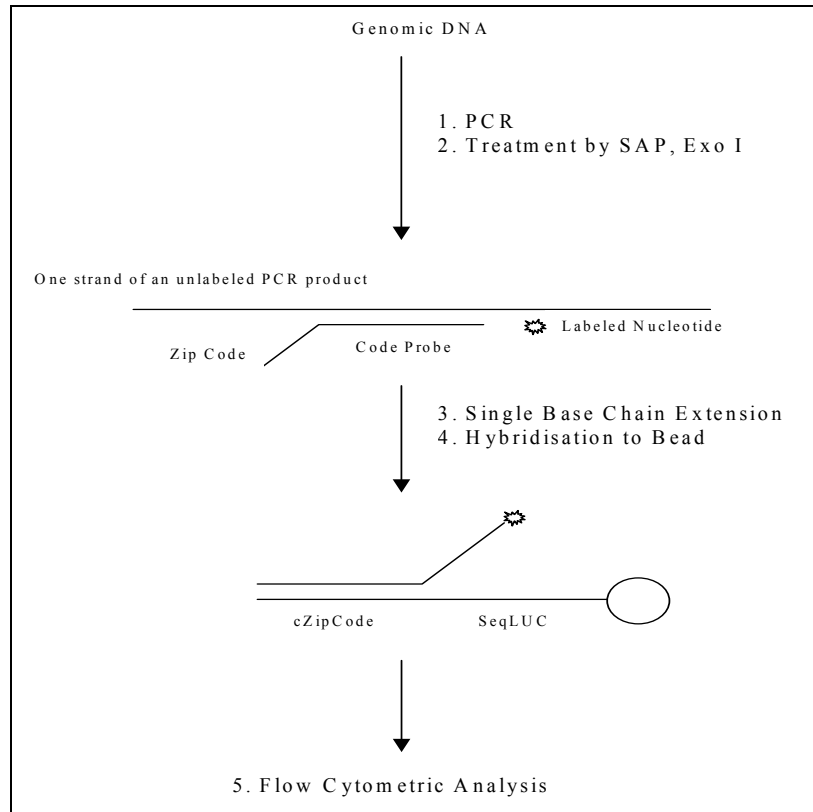


Figure 1.3. Principle of the microsphere single base chain extension assay (SBCE). Firstly, a PCR is done. Then the reaction mix is treated with shrimp alkaline phosphatase (SAP) and exonuclease I (Exo I). In the third step the allele-specific products are generated by primer extension with different labelled fluorescent ddNTPs. For every SNP, a primer with a unique ZipCode sequence is designed and used to detect the two alleles. Multiplex SNP analysis could be done by the use of different ZipCodes for different SNPs employing pooled PCR products. After the extension reaction microspheres are added to the completed SBCE reaction. To these microspheres oligonucleotides encoding the complement to the respective ZipCode sequences (cZipCode) and a common luciferase sequence (SeqLUC) are attached. ZipCode and cZipCode sequences are specifically hybridised on the microspheres and respective fluorescent signals are detected in a flow cytometer.

Each of the microspheres contains a fluorescent colour code. Two unique fluorescent dyes are combined at ten different concentrations providing a set of 100 distinguishable entities. As the pool of microspheres can be rearranged there is a certain degree of flexibility. For SNP genotyping a sequence containing the SNP to be analysed and labelled with a third dye is hybridised to the complementary sequence attached on the bead. As in the case of microarrays primer extension (figure 1.3) and oligonucleotide

ligation are used to improve the allele-specificity. Photon counting is applied for the genotype analysis.

In one approach a flow cytometer is used for automatic analysis, which can be done in a few seconds. In a further variation the coded microspheres are captured in solid phase wells, which are coupled to a fibre optic detection system.

### *The TaqMan assay, “Molecular Beacons” and Kinetic PCR*

The TaqMan assay (figure 1.4) is based on the principles of the fluorescence resonance energy transfer (FRET) system that requires two different, linear oligonucleotide probes.<sup>50</sup> The 3'-end of one oligonucleotide carries a donor fluorophore while the 5'-end of an adjacent oligonucleotide carries an acceptor fluorophore. During FRET, the donor emits photons that are absorbed by the acceptor that then emits fluorescence. FRET is only possible if the distance between donor and acceptor is no more than 6 nucleobases. All of the methods described in this and the following chapter are performed in a homogenous format.

TaqMan probe oligonucleotides anneal between the upstream and the downstream primer in PCR.<sup>51</sup> They contain a fluorophore at the 5'-ends and a quencher at the 3'-ends. As long as the fluorophore and the quencher are linked to the oligonucleotide, the fluorescence is quenched. During PCR amplification with a DNA-polymerase containing 5'-exonuclease activity, the fluorescently labelled 5'-terminal base of an overhang of a probe is cleaved off in the case of complete complementarity. The cleavage suppresses fluorescence quenching, thus leading to fluorescent light. The quantity of fluorescence is directly proportional to the amount of the accumulating PCR product.

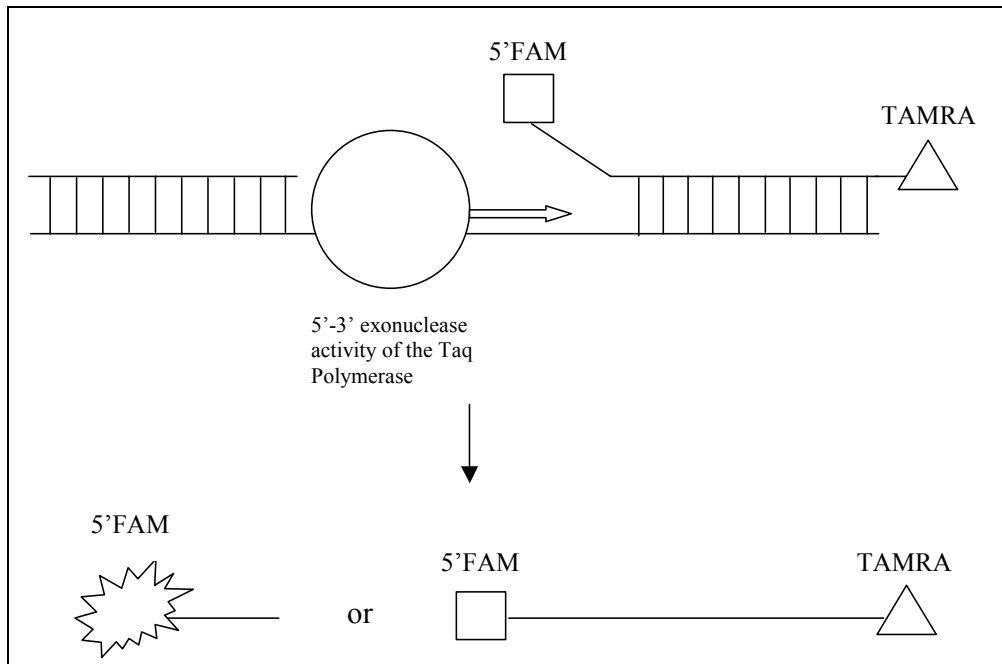


Figure 1.4. Principle of the TaqMan assay. Details are described in the text.

A related assay is termed “Molecular Beacons”. In contrast to TaqMan probes they contain a stem-loop structure, a fluorophore and a quencher (figure 1.5).<sup>52,53</sup> These oligomers are added as detectors to a PCR.

The stem sequence of a molecular beacon is unrelated to the DNA target and only keeps the fluorophore and the quencher in close proximity, while the loop sequence is complementary with the target sequence.

If the probe finds its target, the loop opens and one of its sequences hybridises to the target. Therefore the fluorophore is removed from the vicinity of the quencher releasing fluorescent light. Due to its stem-loop structure the molecular beacons show a higher specificity than linear oligonucleotide probes.

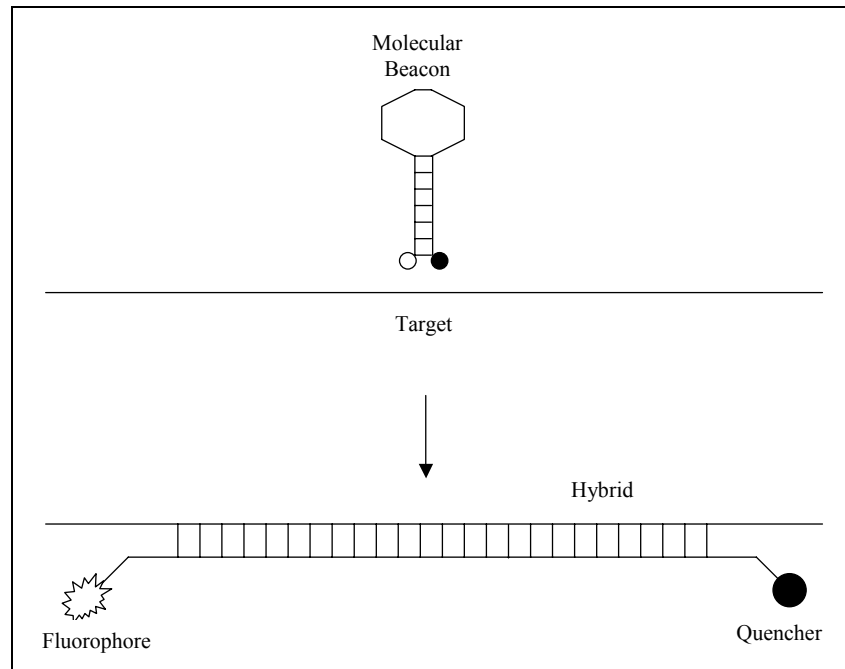


Figure 1.5. Principle of the “Molecular Beacons”. More details are described in the text.

The disadvantages of these systems seem to be the limited degree of detection channels, a demanding automation and a too high price of 1 US-Dollar per SNP analysis.<sup>54</sup> Certain DNA sequences were difficult to differentiate due to similar affinity of perfect-match and mismatch probes for the target sequence.

Another fluorescence-based approach for SNP genotyping is allele-frequency determination in pooled DNA samples by kinetic PCR.<sup>55</sup> This method combines real-time quantitative PCR with allele-specific amplification. Placing the 3'-end of one of the two allele-specific primers directly over the SNP position and matching one of the alternative nucleotides ensures the specificity of the PCR reaction. Allele-specific amplification is observed by increasing fluorescence of DNA-binding dyes like SYBR Green I. Generally, mismatch amplification is delayed by more than 10 cycles. As for kinetic PCR samples have to be analysed after every PCR cycle it is a time-consuming procedure but reagents are cheaper than for TaqMan or Molecular Beacons.

*The Invader assay*

The Invader assay is done in a homogenous format. It is based on the unique capability of a class of natural enzymes called flap-endonucleases and engineered enzymes termed cleavases. DNA molecules are cleaved at specific structures produced by the addition of certain oligonucleotides to DNA or RNA.<sup>56,57</sup>

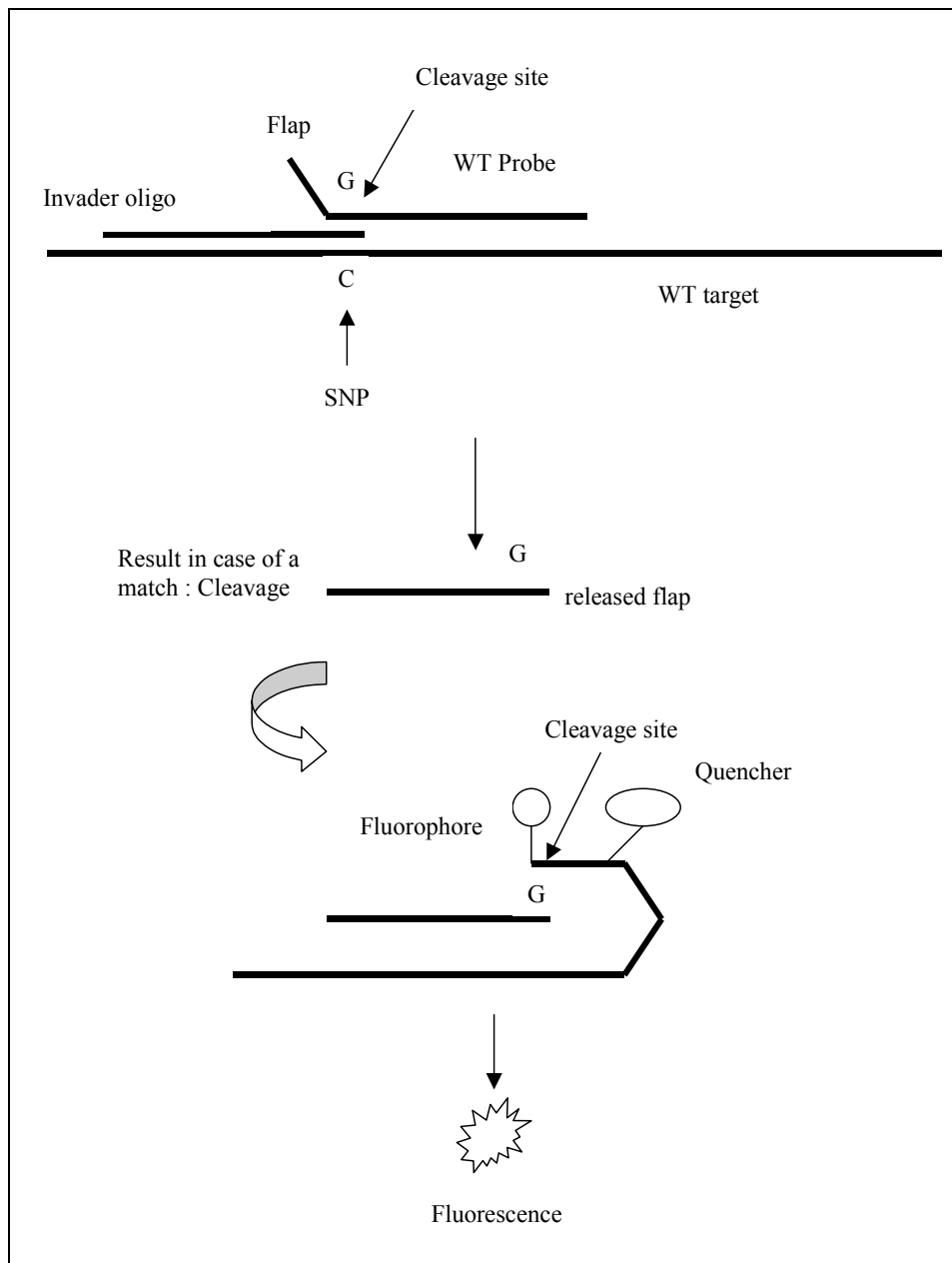


Figure 1.6. Principle of the Invader assay. Details are described in the text.

As is displayed in figure 1.6, the reaction works in such a way that two oligonucleotides hybridise in an adjacent manner to the target nucleic acid.<sup>58</sup>

The upstream oligonucleotide is called invader, while the downstream oligonucleotide is called probe. The resulting adjacent duplexes have to overlap by at least one nucleotide to create an efficient substrate. The 5'-end of the probe contains an unpaired region at the 5'-end called the "flap". The flap will be released by cleavage as a target-specific product if the correct structure was formed before. Specific cleavage of the probe occurs at the position defined by the 3'-end of the upstream oligonucleotides that displaces or "invades" the probe. If the overlap between invader and probe is only one nucleotide, cleavage between the first two base pairs at the 5'-end of the probe takes place, thus releasing the flap and one nucleotide of the base pair region. If the correct invasive configuration is not formed, for example in the case of a mutant DNA target with a wildtype probe, cleavage will not occur. The released flap serves in a subsequent step as an invader oligonucleotide on a FRET probe that is 5'-end labelled with fluorescein and quenched by an internal dye. Hence the procedure is termed "squared invader". Upon cleavage the 5'-fluorescein labelled product is detectable.

The reactions are executed close to the melting temperatures of the probes. An advantage of the Invader assay is that it does not require thermocycling because of its isothermic, balanced equilibrium. Each target-specific product enables the cleavage of many FRET probes. Under standard conditions ca.  $10^6$ - $10^7$  labelled cleaved flaps are produced per hour. Another advantage of this assay is that it works with genomic DNA and does not require PCR amplification thereby avoiding potential contamination problems. During the first invasive cleavage the genomic DNA is the limiting component, since the Invader and probe oligonucleotides are supplied in molar excess. In the second step, the limiting component is the released flap. Several kits using squared invader technology for mutation screening are commercially available but a high-throughput platform for SNP genotyping is not yet established.<sup>59</sup> The Invader assay seems to be difficult to optimise. In particular allele-specific generation of products by the flap-endonucleases depends strongly on buffer conditions, temperatures and target/probe sequences. Probably because of these reasons in large-scale projects applying Invader technology, PCRs were used for the generation of a sufficient amount of template.<sup>60</sup> PCRs were apparently easier to optimise and to multiplex than linear

amplification by the flap-endonuclease.<sup>61</sup> Moreover detailed PCR protocols were accessible. After PCR amplification the allele-specific reaction was carried out by a flap-endonuclease.

### *Introduction to MALDI-MS*

In principle, mass spectrometry provides one of the most attractive solutions for SNP genotyping because it can be used to obtain direct and rapid measurement of DNA. Therefore it is very popular on the detection front as its results can be scored easily and rapidly by automated data management systems. Particularly matrix-assisted laser desorption/ionisation time-of-flight mass spectrometry (MALDI-MS) has revolutionised the mass spectrometric analysis of biomolecules.<sup>62</sup> It was discovered that irradiation of crystals formed by suitable small and acidic organic molecules, termed the “matrix”, with a pulsed laser at a wavelength close to a resonant absorption band of the matrix induced an energy transfer and desorption process, evaporating matrix ions into the gas-phase. The matrix must absorb at the wavelength of the laser - generally a nitrogen laser - for ionisation to occur. Furthermore, it was found that by incorporation of large analyte molecules like proteins into the crystalline structure of the matrix, the non-absorbing molecules were co-desorbed into the gas-phase by fast heating and ionised upon irradiation with the laser. The ionisation process in MALDI is not yet well understood. Several mechanisms for ionisation of large molecules were suggested.<sup>63</sup> By MALDI predominantly either positive or negative single charged molecules are detected.<sup>64</sup> These ions are produced by a proton-transfer reaction of matrix and analyte molecules in the gas-phase. The ions are accelerated by an electric field. Usually, MALDI-MS is performed with time-of-flight separation (MALDI-TOF-MS).<sup>65</sup> Molecules are guided by ion optics into a flight-tube where they are separated before they finally reach the detector.

MALDI-MS has been applied in different variations for the analysis of proteins, peptides and nucleic acids.<sup>66</sup> Its main advantage over conventional DNA diagnostic methods is its speed of signal acquisition (around 100 microseconds for one signal) and



its accuracy of the signal because the signal obtained is the molecular weight, a physical and intrinsic property. In contrast to this, conventional electrophoretic methods for separating and detecting DNA take hours to complete. Additionally, these methods and hybridisation techniques such as microarrays are susceptible to complications deriving from secondary-structure formation in nucleic acids. A further advantage is that no fluorescent dyes, which are expensive, are required. As the complete automation from sample preparation to the acquisition and processing of data is possible, MALDI-MS is generally considered to be an ideal analysis method for high-throughput applications like SNP genotyping.<sup>67</sup>

The principle construction of a MALDI-TOF mass spectrometer is shown in figure 1.7. The equipment of MALDI-TOF-MS instruments with a delayed extraction has greatly improved the resolution of MALDI signals.<sup>68,69</sup>

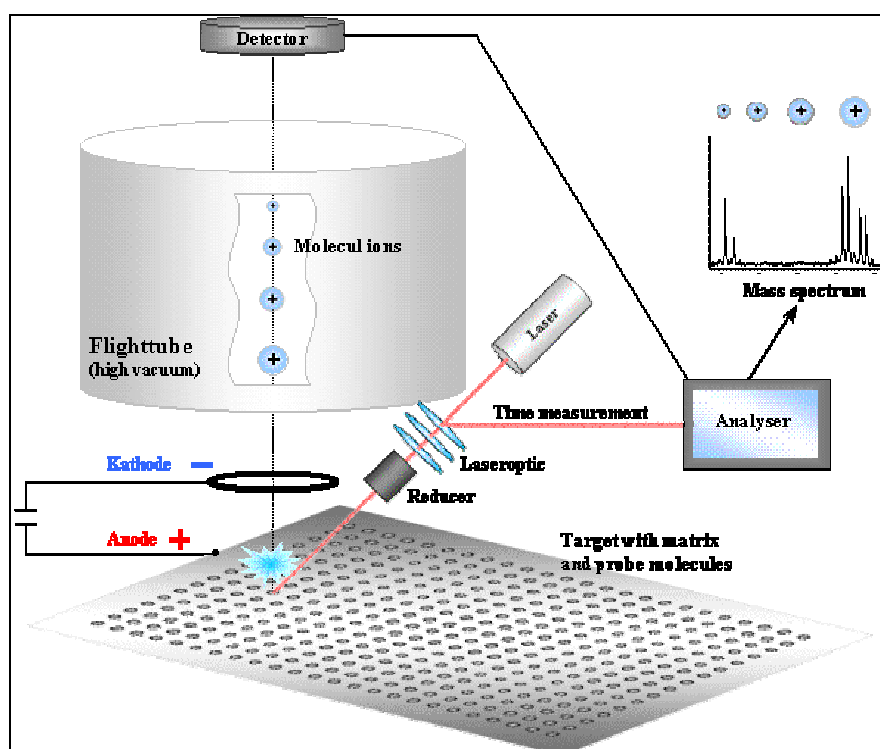


Figure 1.7. Principle of MALDI-TOF-MS. The MALDI features are simplified and the scale is not proportional. More details are described in the text. This picture was kindly provided by Ole Brandt.

A further advance in instrumentation could be the introduction of an IR- instead of the commonly used UV-laser.<sup>70</sup> This is a less aggressive laser therefore desorbing matrix and analyte molecules more carefully. By this enhancement in instrumentation the measurement of DNA with a size up to 1,000 bases was shown. Unfortunately, routine with IR mass spectrometers is difficult and therefore these mass spectrometers are currently not commercially available.

The different types of DNA analysis by MALDI-MS range from the analysis of PCR products to procedures using allele specific termination to single nucleotide primer extension reactions, hybridisation and sequencing.<sup>71</sup> All of these approaches that are described in more detail below heavily rely on stringent purification procedures prior to MALDI analysis. Spin column purification, magnetic bead technology or reverse-phase binding is applied which is cumbersome to use, expensive, with significant batch to batch variation and therefore not easy to incorporate into an automated high throughput set up.<sup>28</sup>

### *DNA sequencing for SNP genotyping by MALDI-MS*

After its invention MALDI-MS was proposed as an alternative for gel-based analysis of DNA sequencing products.<sup>72</sup> It was conceivable that this approach could be used to genotype SNPs. Indeed detection of DNA sequencing ladders by MALDI-MS was demonstrated.<sup>73,74</sup> However, several studies revealed a loss of signal intensity and mass resolution with increasing DNA size.<sup>75,76</sup> Because of the size-dependent loss of signal MALDI-MS is limited to DNA molecules smaller than 100 nucleotides.<sup>77</sup> These disadvantages significantly limit the power of DNA sequencing by MALDI-MS. One contributing factor to this might be the size-dependent tendency of the phosphodiester backbone of the DNA to fragment during the MALDI process, which results in a loss of signal intensity for intact DNA.<sup>78</sup> An additional contributing factor to this restriction could be a bias of MALDI-MS towards smaller DNA molecules. Also the increased sodium and potassium adduct formation of larger DNA fragments, which results in division of their signal over several peaks and an ionisation bias favouring the ionisation

of smaller oligonucleotides causes mentioned problems.<sup>63</sup> To counteract this, stringent purification was performed. Primers for sequencing reactions were employed containing a biotin group that binds to streptavidin-coated magnetic beads required for separation of the reaction products.<sup>79</sup>

### *Direct mass-analysis of PCR products by MALDI-MS*

The described limitations of MALDI-MS have also complicated the detection of PCR amplicons containing SNPs.<sup>80</sup> The successful analysis of a single-stranded amplicon with a size of 69 bases containing one SNP has been shown. Nevertheless, an experiment like this remains difficult to perform, particularly for routine high-throughput analysis. As double-stranded PCR products generally dissociate during the MALDI process into single strands of slightly different masses the resulting signals are poorly resolved and peak broadening and mass inaccuracies are the rule.<sup>81</sup> Masses as small as 9 Da (the mass difference between Thymin and Adenin) are impossible to resolve at 30,000 Da. One way to circumvent these problems was to analyse DNA stretches that were produced in allele-specific PCRs.<sup>82</sup> Primers of these PCRs were constructed to be of sufficiently different masses for easy peak distinction in a mass spectrum. Analogous to the procedure for DNA sequencing stringent purification of the PCR products is essential.

### *Primer Extension and MALDI-MS*

The following procedures for SNP analysis use a primer extension reaction to generate allele-specific products. The advantage of this strategy is that product masses of around 5,000-6,000 Da are definitely smaller than those of the two preceding approaches. A primer is chosen upstream of the SNP that is to be genotyped. A reaction of a PCR amplicon with an extension primer, dNTPs and/or ddNTPs, and a DNA polymerase results in allele specific primer extension products for MALDI detection.

The polymerase extends the 3'-end of the primer by specifically incorporating nucleotides that are complementary to the DNA template. Extension terminates at the first base in the template where a nucleotide occurs that is complementary to one of the ddNTPs in the reaction mix. Some protocols use primers or ddNTPs containing mass-tags that increase the mass differences between the products.<sup>83,84</sup> Generally, a thermostable DNA polymerase in a temperature-cycled reaction is employed, which leads to a linear amplification of extended primers. The analysis of a synthetic oligonucleotide template with a concentration as low as 400 pM was shown.<sup>85</sup> However, the primer extension on a double-stranded PCR product is usually performed at concentrations in the micromolar range. Unfortunately, although several primer extension protocols for MALDI-MS have been developed all of them require stringent solid-phase purification, which contributes to the major cost for high-throughput SNP genotyping and is cumbersome for automation.<sup>86-90</sup>

An interesting variation integrates the primer extension procedure into a semi-automated system called "MALDI on a chip technology" or "MassArray" making use of piezopipetting.<sup>91</sup> Sample preparation is achieved using microdispenser nozzles, which deliver droplets by a pulsed voltage over a piezo-ceramic element. Only some nanoliters of a sample from the molecular biological reaction are pipetted onto a silicon chip that is inserted directly into the MALDI mass spectrometer where each sample spot is measured automatically.

Currently, primer extension has become the most widely used molecular biological procedure for SNP analysis by MALDI-MS because of its allele-specificity and generation of fairly small products, which also holds for alternative detection methods. Well known commercially available systems for SNP genotyping using the strategies presented here are provided by Sequenom (San Diego, CA) and Perseptive (Framingham, MA).<sup>92,93</sup>

*Nucleic acid hybridisation and mass spectrometry*

Peptide nucleic acid (PNA) is a DNA analogue containing a charge neutral amide backbone and the four regular nucleobases (figure 1.8).<sup>94,95</sup> Because of the modified backbone, PNA is not degraded by nucleases and therefore it might be useful for antisense applications such as expression regulation. Additionally, the amide backbone of the PNA has several advantages for allele-specific hybridisation compared to unmodified DNA, for example an increased thermal stability of the duplex (PNA/DNA), the ability to hybridise under low ionic strength conditions and higher hybridisation specificity for complementary DNA probes.<sup>96</sup> Furthermore, PNA is more easily analysed by MALDI-MS than DNA. The PNA backbone, in contrast to DNA containing a negative charge backbone, does not fragment easily during the MALDI process and does not tend to form metal ion adducts.

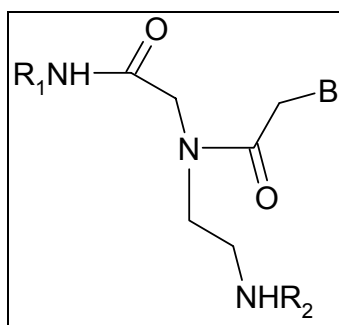


Figure 1.8. Chemical structure of PNA

Two quite similar procedures using PNA hybridisation have been developed for SNP genotyping by MALDI-MS.<sup>97,98</sup> In both cases, two PNA probes, each corresponding to one allele, were hybridised to biotinylated, single-stranded PCR products that were immobilised on streptavidin-coated magnetic beads. Afterwards the beads were washed stringently so that only the perfectly matched PNA remained annealed to the PCR product. Then the bead solution was spotted onto a MALDI probe tip and a matrix solution, which is acidic, was added to the beads in order to dissociate the PNA/DNA duplex. The PNAs were successfully desorbed and ionised with the matrix upon laser irradiation and analysed by MALDI-MS. In one variation, each allele-specific PNA

probe was mass-tagged by incorporating variable numbers of 8-amino-3,6-dioxaoctanoic acid residues on the N-terminal end.<sup>97</sup> In another approach, mass tagging of the PNAs was done by simple adding of an extra, non-complementary dT nucleobase to the 3'-end of the PNA.<sup>98</sup> Unfortunately, PNAs have a high variability concerning thermal stabilities of different sequences. This requires respective optimisation of the hybridisation of probe pairs for each SNP. Nevertheless, it was claimed that by the use of different algorithms for the prediction of PNA/DNA duplex stabilities, single tube multiplex SNP genotyping would be possible.<sup>99,100</sup> Additional disadvantages are the extremely high price of PNA components, which generally cost 10-fold more than analogous DNA compounds, and the use of expensive magnetic beads.

A completely different approach applies hybridisation of oligonucleotide probes representing genotypes on purified PCR products and nuclease digestion for selection.<sup>101</sup> A complementary, tight-bounded probe will resist the nuclease digestion, while phosphodiesterase I attacks non-complementary probes. Under the employed matrix conditions the surviving probe is detected. A current disadvantage of this method might be that a lot of laser shots (4 times 400!) had to be accumulated. A more principle problem might lie in the specificity of the procedure.

An alternative approach using DNA hybridisation is used by Masscode technology offered by Qiagen.<sup>102</sup> Instead of a MALDI mass spectrometer an electrospray ionisation quadrupole mass spectrometer (ESI-MS) is employed. The use of DNA polymerases in allele-specific PCRs, which highly depend on hybridisation, increases the generation of allele-specific products. Thus if the primers hybridise perfectly amplification occurs. Enzymatic amplification has two advantages. It increases allele-specificity and generates enough products for mass spectrometric detection. Primers used for allele-specific PCR consist of different 5'-mass-tags and a photocleavable linker. Currently more than 100 different mass tags have been synthesized. Each tag serves as a unique discriminator for an oligonucleotide used for allele-specific PCR. After photo-cleavage by ultra-violet light the respective mass tag is analysed in the mass spectrometer. The price per SNP analysis is currently in the range of 1 US-Dollar.

### *Invader assay by MALDI-MS*

As explained above, the enzymatic procedure involves the sequence specific hybridisation of two oligonucleotides to form an overlapping structure at the SNP to be studied. In the case of hybridisation an enzymatic cleavage with a thermostable flap endonuclease is followed and an allele-specific, short oligonucleotide signal molecule deriving from its overlap-structure, the “flap”, is linearly amplified. In contrast to the fluorescence-based Invader assay the flaps are analysed directly by their specific masses. The 5'-ends of the flaps contain a biotin group allowing binding on streptavidin-coated magnetic beads.<sup>103,104</sup> After purification, clean DNA probes are eluted directly for MALDI sample preparation and allele-specific flap oligomers are then analysed in the mass spectrometer.

### *Problems of (UV-) MALDI-MS analysis*

MALDI was initially applied to the analysis of peptides and proteins. In contrast to this, DNA is significantly more difficult to analyse because of its chemical structure and properties.<sup>78</sup> The main problem in analysing native DNA by MALDI-MS consists in its negatively charged sugar-phosphate backbone. With native DNA, the phosphate residue provides a site of negative charge in solution and each DNA molecule carries as many negative charges as phosphate residues. The affinity of the phosphate residues for alkali counterions, such as sodium and potassium, or other metal counterions, is high, but not high enough to result in complete saturation. These ions interfere with the ionisation process, by inducing adducts and thereby significantly reducing the signal intensity.<sup>105</sup> The use of ammonium counterions in MALDI is a well-established method to counteract ion affinities.<sup>106</sup> In solution ammonium exists as a  $\text{NH}_4^+$  counterion, whereas in the gas-phase  $\text{NH}_3$  is readily lost, leading to a reduced counterion structure. However, ammonium ions introduce a degree of suppression to the desorption process. This results in a decrease of analytical resolution and sensitivity. Nowadays stringent purification procedures are applied to overcome these problems including magnetic bead separation and reversed-phase column purification.

Another counteracting feature is the acid instability of DNA. Sample preparation is performed with acidic matrices but acidic conditions are encountered in the desorption/ionisation process. In the gas-phase, DNA can readily fragment with harsh matrices. A detectable degree of depurination has been observed for larger DNA products.<sup>107</sup> Replacing purines by 7-deaza-analogues is one approach to prevent DNA from depurinating.<sup>108,109</sup> A second procedure to improve DNA in MALDI is the use of ribonucleotides containing 2'-OH groups that stabilise the gas-phase ion.<sup>110</sup> In a third approach it was found that the replacement of phosphate protons from native DNA backbones by alkyl groups significantly improved the behaviour of the molecule in the MALDI process.<sup>111,112</sup>

The optimisation of the MALDI process consists of identifying the right matrix and preparation method for an analyte. How matrices function in MALDI is until now not well understood. The chemical structure of DNA is complex and its interaction with a matrix during the desorption/ionisation process eludes investigation. Only empirical findings progressed the method. It was observed that DNA analysis by MALDI was very inefficient, for example 100 times more DNA has to be used in a preparation to achieve similar signal intensities comparable to peptides.<sup>111</sup>

The principal idea for rendering DNA amenable to analysis by MALDI focused on the difference in analysing oligonucleotides and peptides. While most peptides are formally uncharged, DNA carries as many negative charges as phosphate bridges. Charges were neutralised by replacing phosphate groups by phosphorothioate groups and alkylating them. The efficiency of alkylation of regular phosphate groups is low, but a selective and quantitative alkylation is achieved with phosphorothioate groups. Furthermore, it was known that the addition of a positive charge tag to peptides improved their desorption behaviour.<sup>113</sup> Therefore the addition of a positive charge tag with subsequent removal of all charges from the phosphorothioate backbone bridges was implemented to the analysis of DNA by MALDI-MS. This procedure was called "charge-tagging". The concept of this was to generate a product with a defined charge state, thus relying on the matrix for desorption, but not for ionisation. Using this approach, there was a 100-fold increase of detection efficiency, equalling the detection efficiency of peptides.<sup>114</sup> The same result was observed when all but one backbone-



bridge was neutralised and the DNA product thus carried a single negative charge (-1 charged DNA product).<sup>115</sup>

$\alpha$ -Cyano-4-hydroxy-cinnamic acid methyl ester was found to be the ideal matrix system for DNA compounds with either one single positive or one single negative charge. It is the methyl ester of  $\alpha$ -cyano-4-hydroxy-cinnamic acid, the most commonly used matrix for peptide analysis. In contrast to other matrices it has a significantly higher  $pK_a$  of around 8 in solution. Its absorption maximum perfectly matches the emission wavelength of an  $N_2$  laser, which is the most commonly used laser in MALDI mass spectrometers. In contrast, matrices used for protein and peptide analysis typically have a very low  $pK_a$ . Standard DNA matrices, like 3-hydroxypicolinic acid (HPA, the most common matrix for DNA analysis) have slightly acidic  $pK_a$ 's in the region of around 4. One of the most striking observations with  $\alpha$ -cyano-4-hydroxy-cinnamic acid methyl ester was that native DNA could not be analysed with this matrix.<sup>114</sup> Use can be made of this discriminative behaviour as the selectivity of this matrix is towards singly charged DNA compounds. It was claimed that there is little difference in ionisation efficiency in negative or positive ion mode analysis of singly charged oligonucleotides with this matrix.<sup>115</sup>

Two common matrix preparation methods are applied in MALDI-MS, thin-layer and dried droplet preparation. For thin-layer preparations the matrix is spread over the MALDI target plate in a volatile solvent, such as acetone. The solvent evaporates immediately, leaving a thin layer of small matrix crystals. The analyte is dispensed onto the thin-layer in a solvent that does not dissolve the matrix. Analyte molecules co-crystallise into the surface of the matrix. Hence the analytes are desorbed approximately equally all over the spot leading to better mass accuracy. For dried droplet preparations a matrix solution is mixed with an analyte solution and then spotted onto the MALDI target plate. Dried droplet preparations result in "sweet spots". Certain positions on the preparation give better results than others, which make these preparations quite difficult to use in automated processes. Due to the uneven height of dried droplet preparations the mass calibration can also be unstable. MALDI analysis is based on the determination of the time-of-flight of an ion. Variable height of the matrix preparation results in a shift of the starting position what affects the time of flight. This can easily conclude in a few Daltons mass variation. In contrast, thin layer preparations give less

spot-to-spot variation, better mass accuracy and resolution. Thin-layer is used with  $\alpha$ -cyano-4-hydroxy-cinnamic acid for peptide analysis, while DNA analysis preferably is done with HPA in a dried droplet preparation.

### *Objectives of this thesis*

An ideal method for SNP genotyping would be homogenous, easy in handling, efficient in reagent consumption, highly specific in readout, and the results interpretable by computer software. MALDI-MS is considered to be a very powerful technique for DNA analysis because of its speed and accuracy.<sup>28,29,71</sup> Nevertheless, MALDI-MS of DNA requires stringent purification procedures, which is a big disadvantage for automation and contributes to the major cost of the analysis. A procedure that uses MALDI-MS combined with an automation process that does not require any purification steps is therefore very sought-after.

The know-how of sensitivity enhancing chemistry for DNA analysis by MALDI could help to cope with mentioned problems. There exist no methods to produce allele-specific DNA molecules that introduce the described DNA modifications. Preferably this could be done enzymatically, as enzymes such as DNA polymerases, allele-specific endonucleases or ligases could provide high specificity. The interface of molecular biology and DNA modifications comprising charge-neutral DNA backbones and “charge-tags” should be studied. Based on this purification-free procedures should be established. The stability for high-throughput, multiplexibility and variability for easy optimisation of daily SNP assay development should be evaluated. The final objective is that the developed procedures terminate in efficient automated processes for high-throughput SNP genotyping.