

Chapter 1

Introduction

In this chapter we examine the role and the importance of digital certificates in communication and transaction mechanisms. We discuss the main developments and point out their security, efficiency, and privacy shortcomings. Next we examine the meager previous efforts to protect privacy in public key infrastructures. Amongst others, we show that the popular suggestion to offer privacy by issuing pseudonymous certificates is not only insecure in almost all situations, but also ineffective to protect privacy. On the basis of the previous findings we list basic desirable privacy properties. Finally, we outline how the techniques that will be developed in later chapters meet these and other privacy properties and at the same time help overcome the security and efficiency problems.

1.1 Digital certificates and PKIs

1.1.1 From paper-based to digital certificates

Individuals and organizations often have a legitimate need to verify the identity or other attributes of the individuals they communicate or transact with. The traditional method for demonstrating that one meets certain qualifications is to disclose one or more paper-based certificates. As defined in the third edition of the American Heritage Dictionary of the English Language, a certificate is “a document testifying to the truth of something.” Photographs, handwritten signatures, and physical cues help the verifier to establish the identity of the holder of a certificate. Embedded security features (such as special paper, watermarks, ink that appears different when viewed from different angles, and microprinted words and other detail that is hard to replicate) serve to protect against counterfeiting and unauthorized duplication.

Since the advent of computers and telecommunication networks, paper-based transaction mechanisms are being replaced by electronic transaction mechanisms at

a breath-taking pace. Many forces drive this unstoppable transition:

- The theft, loss, or destruction of a paper-based certificate coincides with the theft, loss, or destruction of at least part of its value. It may be expensive, difficult, or impossible to obtain a new copy from the issuer.
- Paper-based certificates are subject to wear and tear, add to the depletion of forests, are costly to handle, and in many situations are inefficient. Electronic certificates can be manufactured, distributed, copied, verified, and processed much more efficiently and at lesser cost.
- Paper-based certificates are not suitable to convey negative qualifications of their holders. An individual carrying a certificate attesting to the fact that he or she has been in prison, say, can simply discard the certificate. Sometimes negative qualifications can be tied in with positive ones (e.g., a mark for drunk driving on a driver's license), but this measure is not always an option.
- Cyberspace (the conglomeration of networks that enable remote communication, including the Internet, e-mail, cable TV, and mobile phone networks such as GSM) offers huge benefits over face-to-face communications and transactions in the physical world. Many of the benefits cannot be realized using paper-based certificates, however, since these require physical transport.
- The public at large can avail itself at modest cost of ever-advancing desktop reprographic equipment. A nationwide study conducted in 1998 by U.S. corporate investigation firm Kessler & Associates found resume and credential fraud to be of "almost epidemic proportions." Counterfeiting rarely requires perfection; it usually suffices to produce something that will pass casual human inspection. Ultimately, the counterfeiting threat can be overcome only by moving to certificates that are cryptographically secured and that can be verified with 100 percent accuracy by computers.

In many applications, symmetric cryptographic techniques are inappropriate: they require a trusted third party to set up a secret key for any two parties that have not communicated previously, and cannot offer non-repudiation. Thus, there is a fundamental need for public key cryptography. Public key cryptography enables the parties in a system to digitally sign and encrypt their messages. When two parties that have not communicated before want to establish an authenticated session, they need merely fetch the public key of the other; there is no need for a trusted third party to mediate every transaction.

In their seminal paper [136] on public key cryptography, Diffie and Hellman pointed out the problem of authenticating that a public key belongs to an entity. They suggested using secure online repositories with entries that specify name-key bindings. In 1978, Kohnfelder [238] proposed to avoid this potential bottleneck by having

a trusted entity, called the Certificate Authority¹ (CA), vouch for the binding between a public key and its holder. A *digital certificate* is a signed assertion about a public key. More specifically, it is a digital signature of the CA that binds a public key to some other piece of information, in Kohnfelder's case the name of the legitimate holder of the public key. This enables all system participants to verify the name–key binding of any presented certificate by applying the public key of the CA. There is no need to involve the CA in the verification process; verification can be off-line.

A *public key infrastructure* (PKI), also called key management infrastructure, is an infrastructure for a distributed environment that centers around the distribution and management of public keys and digital certificates. It is widely recognized that PKIs are an essential ingredient for secure electronic communications and transactions in open environments. See, for instance, Feghhi, Williams, and Feghhi [167], Ford and Baum [172], Froomkin [177], Lewis [251], and Zimits and Montano [395].

The CA can be made responsible not only for certifying public keys and authenticating certificate applicants, but also for notarizing electronic documents, resolving disputes, and keeping track of revoked keys. Some or all of these functions may be managed by separate trusted parties. For instance, the registration and approval of certificate applicants may be done by a separate Registration Authority.

In practice, a PKI can have multiple CAs, so that certificate applicants and verifiers need not trust a single CA. CAs can certify the public keys of other CAs, and in this manner arbitrary CA structures can be formed. This gives rise to such notions as certificate chains, bridge CA's, and cross certification; see Burr [68]. Our techniques enable anyone to be the issuer of their own digital certificates, and all issuers can coexist in a single PKI. We will not address multi-CA PKIs, though, because the techniques for these are straightforward and largely orthogonal to the techniques that we will develop in this book. For simplicity, we will always assume that each PKI has only a single CA, unless explicitly stated otherwise.

Also for simplicity, we will often equate certificate holders with individuals, and certificate verifiers with organizations. More generally, the entities that retrieve, hold, show, verify, or otherwise operate on certificates may be software programs, hardware devices, or anything else that can perform the required logical steps.

1.1.2 Identity certificates

An object *identifier* is any data string that can readily and uniquely be associated with the object. What Kohnfelder called a digital certificate is better referred to as an *identity certificate*, because it binds a public key to a person identifier, such as a credit card number, a “true name,” a fingerprint, a Social Security number, or a health registration number.

The X.509 certificate framework [216] is the best known example of identity certificates. In 1988, the International Telecommunications Union (formerly the In-

¹In recent years the term Trusted Third Party (TTP) has gained in popularity.

ternational Consultative Committee on Telephone and Telegraphy) started working on X.509. X.509v1 was designed to certify the public keys of principals that are uniquely named in X.500 [80, 195, 385], an online database listing globally unique names; an entry in an X.500 directory can be a person, a device, or anything else that can be assigned a “Distinguished Name.” X.509v2, released in 1993, provided for a more flexible choice of identifiers. X.509v3, announced in June 1997 (see [218] for amendments), greatly improved the flexibility of X.509 certificates, by providing for a generic mechanism to extend certificates. Also, X.509v3 allows the use of local names in certificates, acknowledging that a global naming scheme is unworkable.

Numerous (draft) standards and CA products have been developed based on the X.509 framework. X9.55 [8, 11], for example, is an ANSI-adopted standard developed by the American Bankers Association that is similar to X.509 but targeted at the financial services industry. Another effort is PKIX [4, 114], a draft standard by the Internet Engineering Task Force (IETF) to make X.509v3 certificates suitable for the Internet.² Other implementations of X.509 certificates include Privacy Enhanced Mail [86, 230] (PEM, an IETF e-mail standard proposal), Fortezza (the standard for secure e-mail and file encryption in the U.S. defense system), Secure/Multipurpose Internet Mail Extensions [140, 141] (S/MIME, an e-mail standard proposed by RSA Security), Secure Socket Layer version 3.0 [174] (SSL, developed by Netscape to support server and client authentication and session encryption), and Secure Electronic Transactions [257] (SET, proposed by MasterCard and Visa for securing card-not-present credit card transactions).

Also, virtually all the pilot PKI projects conducted by 24 U.S. federal agencies (including the NSA, the IRS, the FBI, the U.S. Department of Defense, and the Social Security Administration) as part of the Federal Public Key Infrastructure [163] (FPKI) use X.509v3 certificates, with application-dependent extensions. For instance, the U.S. Department of Defense is building a PKI “to ensure the authenticity of digital signatures on contracting documents, travel vouchers, and other forms that obligate taxpayer funds, to authenticate users of information systems, and protect the privacy of transactions over networks;” see the DoD Public Key Infrastructure Program Management Office [138, 139] for details.

Another U.S. federal PKI plan based on X.509v3 certificates is Access Certificates for Electronic Services [377, 378] (ACES), which will provide for public electronic access to government services and information. Furthermore, the Department of Justice, the Department of Defense, the NSA, and NASA formed a government-industry consortium called Security Proof Of Concept Keystone (SPOCK); its goal is to demonstrate commercial and federal PKI solutions in cooperation with security technology providers. According to the National Institute of Standards and Technology (NIST), which is responsible for U.S. federal computer security, the FPKI will be knit together from these and other PKI efforts.

²IBM and its Lotus Development subsidiary in July 1998 started making the source code for their PKIX implementation Jonah available to the public, to promote applications based on PKIX.

Other jurisdictions that are in advanced stages of planning federal PKIs include the United Kingdom (its CLOUD COVER initiative is aimed to stimulate the growth of a government-wide PKI), Australia (the Australian Public Key Authentication Framework, PKAF for short, will result from the Gatekeeper federal infrastructure program and efforts by the Certification Forum of Australia), Canada (in 1995, the Treasury Board endorsed a project called GOC PKI, for Government of Canada Public Key Infrastructure), and Hong Kong (in November 1999, the Hong Kong postal service started issuing identity certificates to most of the 6.5 million residents). All these efforts are compatible with the X.509v3 standard. For a snap-shot overview as of July 1999 of the PKI initiatives in 26 member countries of the Organisation for Economic Co-operation and Development (OECD, an international organization consisting of 29 primarily industrialized countries), see the Working Party on Information Security and Privacy [296] of the OECD.

Dozens of developers around the world specialize in CA products involving identity certificates, most of them based on X.509. Among the major players are VeriSign, Baltimore Technologies, Entrust Technologies, and Thawte Consulting (acquired in February 2000 by VeriSign). In recent years a host of companies joined them in their race to capture the identity certificate market (either products or services), including ABAecom, ActivCard, BelSign, Brokat, Celo Communications, Certco, CertiSign, Chrysalis-ITS, Cryptomathic, GTE CyberTrust, Cylink, Digital Signature Trust Company, Entegrity Solutions, EuroSign, EuroTrust, Frontier Technologies, Gemplus, GlobalSign, Internet Dynamics, Identrus, InterClear, KeyPOST, KeyWitness, Litronic, RSA Security, Sonera SmartTrust, Spyurus, Sun Certificate Authorities, Utimaco, ValiCert, Xcert International, and Zergo. Also, major corporations including American Express, AT&T, Canada Post, CompuSource, Equifax, Hewlett-Packard, IBM, Lotus Development, Microsoft, Motorola, Netscape, and Novell all support the X.509 digital certificate standard.³ To accelerate the adoption of identity certificates, Baltimore Technologies, Entrust Technologies, IBM, Microsoft, and RSA Security in December 1999 founded an alliance that has since been joined by over 40 other companies.

Another well-known scheme based on identity certificates is Pretty Good Privacy [69, 396] (PGP). PGP certificates bind a public key to a common name and an e-mail address. PGP is based on a different *metric of authentication* (see Levien and Aiken [250] and Reiter and Stubblebine [321]) than X.509: anyone in the PGP “Web of Trust” can certify keys.

As these developments show, identity certificates are widely perceived as a fundamental technology for secure electronic communications and transactions. Market surveys confirm this. A study released in March 2000 by the Radicati Group, for instance, estimates that the market for CA software for identity certificates will grow from over 368 million U.S. dollar in revenues by year end 2000 to over 1.5 billion U.S. dollar by 2004. Another survey by IDC expects the market to grow to 1.3 billion

³Most companies offer services and products based on the CA toolkits of a select few.

U.S. dollar in 2003.

Identity certificates will also play a major role in the many plans outside cyberspace to migrate to chipcards. A *chipcard* is a plastic card that has the shape and thickness of a conventional credit card, and that contains one or more embedded integrated circuits. Around the world, public transport organizations, municipalities, health care providers, ministry departments, financial institutions, and other influential organizations are planning to provide all their customers with a chipcard that will be the sole means of participating in their systems. Due to the storage and computation limitations of current chipcard technologies, identity certificates do not yet have a prominent place in many of these plans. However, over time the move towards digital certificates is inevitable, for security reasons; see Section 1.1.6 for details.

1.1.3 Central database paradigm

In many applications with a need for authentication, organizations are not (primarily) interested in the identity of a key holder, but in the confirmation of previous contacts, the affiliation of the key holder to a group, the authenticity of personal data of the key holder, the eligibility or capability of the key holder to perform certain actions, and so on. Identity certificates can be used by organizations as authenticated pointers into central database entries that contain the relevant data, and thus support any such authentication needs. This central database paradigm allows organizations to consult any databases they are interested in, to update database entries as they see fit, and to securely maintain negative data about system participants. It also enables organizations to build profiles of individuals for the purpose of inventory management, direct marketing, and so on.

It is easy to see why the use of identity certificates in conjunction with central database look-up has become the model of choice: until recently, it was expensive or impractical to resort to decentralized computing and distributed databases. The centralized model, however, has many drawbacks for organizations and other certificate verifiers:

- The transaction process requires a sufficient delay to identify and correct frauds and other undesirable conditions. This may result in organizations being unable to serve as many customers as they could otherwise.
- Because certificate holders are not ensured that their transactions will be authorized, significant uncertainty is introduced in the transaction process. Requests may be rejected on the basis of erroneous or irrelevant data, or simply because the online connection fails due to peak load, a natural disaster, or otherwise. (The chances of an off-line terminal failing are much slimmer, and moreover the certificate verifier may take immediate action to overcome the problem.)
- In case the verifying agents of an organization are geographically distributed, central database verification may be expensive (because of telecommunica-

tions cost or the difficulty of dealing with peak load) or simply not an option because of the absence of network connections.

- Requests for central database look-up may be dishonored for any reason and may be expensive (especially if databases are operated by commercial organizations such as consumer reporting bureaus).
- It is increasingly difficult for organizations to protect their online databases against intrusions by hackers and insiders. This exposes organizations to incidents that might incur legal liability or hurt their reputation.
- The trend is for governments to require organizations that handle personal data⁴ to adhere to (legal or self-enforced) privacy standards. Significant compliance costs are involved with personnel training, making databases accessible to external auditors, and so on.
- The possession of data about the personal preferences and lifestyle of individuals enables organizations to discriminate against their customers in all kinds of ways. This increases the scope for false complaints and legislative actions.

It is ironic that digital certificates today are considered by many to be a secure way to provide access to personal data stored in central databases. The practice of looking up data in real time in a central database goes against the philosophy behind digital certificates, which is to allow off-line verification of digital signatures. In many PKIs it is a waste of efficiency to use digital certificates in combination with central database look-up; one might as well do away with digital certificates altogether and simply check the validity of public keys in a central database. Indeed, Wheeler and Wheeler [388] and the Accredited Standards Committee X9 [3] for this reason propose a return to the online key repository model of Diffie and Hellman. (This model cannot protect the privacy of certificate holders, though, as we will see later on.)

The central database paradigm is even less desirable from the perspective of individuals:

- Individuals can be discriminated against on the basis of data that is not relevant for the situation at hand. Such discrimination could go about without the individual being aware of the source of the discrimination, the nature of the data used against him or her, or even the mere fact of the discrimination. A qualified job applicant may be rejected just because some manager who bothered to consult a few databases (such as Internet newsgroup archives) cannot relate to his or her lifestyle. Likewise, individuals soliciting a loan or anyone of a myriad of other services may find their applications turned down because somewhere in the process someone discriminated against them, or in favor of others, on the grounds of irrelevant data.

⁴The OECD defines [294] personal information as “any information relating to an identified or identifiable individual (data subject).”

Material damage may result when personal data is accessed with malicious intent. Stalkers, murderers, and extortioners use address information from credit reports and other sources that reveal consumer data to track down their victims. Blackmailers persuade their victims by threatening to reveal sensitive personal data, and kidnappers and robbers plan when to strike by following the whereabouts of their victims. Many criminals are not concerned about targeting a particular individual, but instead select their victims on the basis of their profile; robbers and blackmailers mainly target wealthy singles, and political aggressors are often interested in individuals with particular political or religious convictions.

- When data records do not reflect an individual's true situation, perfectly eligible individuals may end up losing their insurances, loans, housing, jobs, reputations, and so on. Errors are far from uncommon. For instance, the sixth study of the U.S. Public Interest Research Group [314] on credit report accuracy and privacy issues found that 29% of U.S. credit reports contain serious errors that could result in the denial of credit, loans, or jobs, and that altogether 70% of credit reports contain mistakes. Data in central databases may not reflect an individual's true situation for a number of reasons:
 - A substantial portion of all captured data is outdated. One cannot reasonably expect individuals to inform all database operators each time their personal circumstances change; individuals in developed countries are stored on average in roughly a 1000 databases, most of which are unknown to them.
 - Another portion contains information that was composed by drawing incorrect inferences from other sources of data.
 - Whenever data is conveyed orally or in writing, errors are bound to be made when the data is translated into machine-readable form.
 - Data stored in databases may be modified or destroyed by hackers and other outsiders. With the rise of the Internet, the risks are increasing dramatically. Hackers almost routinely gain access to databases, both commercial and governmental, and are rarely prevented from erasing or modifying data records without leaving a trace. In an infamous hack in the mid 1980s, a hacker broke into the databases of Experian (one of the three largest U.S. credit bureaus) to peak into the credit records of Ronald Reagan, and discovered 63 other requests for Reagan's records, all logged on the same day.
 - Data stored in databases may be modified or destroyed by authorized database users and other insiders. Any organization of substantial size is bound to have employees who are willing to accept bribes or have malicious intentions of their own. A 1998 survey [215] by the Computer Se-

curity Institute found that the attack that was by far the most reported by its respondents (520 security practitioners in U.S. corporations, government agencies, financial institutions, and universities) was unauthorized access by employees.

- Misbehavior by identity thieves often ends up registered in the database entries of their victims. The incidence of *identity fraud* has been rising dramatically since the mid eighties. Since 1996, calls on identity theft have been the number one topic on the hotline of the U.S. Privacy Rights Clearinghouse. For details on identity fraud, see Cavoukian [78], the Federal Trade Commission [164], the General Accounting Office [183], Givens [186], and the U.S. National Fraud Center [392].

Since errors spread throughout the system and accumulate as data is disseminated and merged, victims may find themselves affected by the same errors over and over again.

- Individuals have lost all control over how personal data in databases is becoming available to others. Collectors of personal data are always tempted to sell the data or to provide access to it in other ways (thousands of information resellers already offer their services over the Internet to anyone willing to pay), information brokers and private investigators resort to trickery (“pretexting”) to obtain all kinds of personal data, and most countries around the world have laws that require database maintainers to provide access to law enforcement when presented with a court order or a warrant. Also, personal data increasingly becomes available to others by error. In recent years the popular press has reported on numerous cases whereby commercial organizations (such as providers of free e-mail services, credit bureaus, and Internet merchants) as well as government organizations (including social security administrations, law enforcement, and taxation authorities) inadvertently released sensitive personal data to the wrong parties or to the public at large.

In many cases it is virtually impossible for victims to seek and obtain redress. The basis or source of discrimination, misuse, or other harmful actions may never become known in the first place, and even if it does, it may be very hard to repudiate the action.

1.1.4 Attribute certificates

In the early 1990s, the idea of *attribute certificates* gained interest. An attribute certificate binds a public key to one or more *attributes*, which X.501 [81] (also known as ISO/IEC 9594-2) defines as “information of any type.”⁵

⁵This terminology makes sense when considering the dictionary meaning of “attribute.” The third edition of the American Heritage Dictionary of the English Language defines an attribute as “a quality or characteristic inherent in or ascribed to someone or something.”

Attribute certificates are a generalization of identity certificates (an identifier is just one of infinitely many attributes), and have naturally evolved from them. Indeed, identity certificates typically specify other data than just a person identifier and a public key. For instance, an X.509v3 certificate also specifies a version number, a serial number (for revocation purposes), a signature algorithm identifier, a CA name, a validity period, a subject name, a CA signature algorithm identifier, a subject public key algorithm identifier, and (optional) CA and subject identifiers and extensions. However, identity certificates typically contain no other personal data than a person identifier.

From now on we reserve the term attribute certificates to refer to digital certificates that serve primarily to enable verifiers to establish attributes other than the identity of the key holder (such as access rights, authorities, adherence to standards or legal requirements, privileges, permissions, capabilities, preferences, assets, demographic information, and policy specifications).

Attribute certificates have important advantages over identity certificates:

- It is inconvenient for millions of individuals to make a physical appearance before CAs. In November 1998, market researcher INTECO Corp. found that only 64% of Internet users would be willing to appear in person to have their identity verified for a digital certificate. For many types of attribute certificates, there is no need to show up in person at a CA.
- It is typically much harder, more error-prone, and more costly for a CA to establish a person's identity than to establish authorities and other personal attributes. In PKIs where organizations are interested only in non-identity attributes, not including identities can therefore bring substantial savings in cost and time, and can reduce the risk of identity fraud.
- Identity certification may expose a CA to much greater liability. Typically, only government agencies and major organizations such as credit bureaus and financial institutions are in a good position to take on the role of establishers of identity.⁶ Indeed, Kaufman Winn and Ellison [228] argue that the CA cannot legally make users liable for actions for which they cannot reasonably be expected to control the risks and losses, because PKIs cannot subsume risk at the technical level. (The latter observation is also at the heart of critiques against

⁶Recent market developments are in line with this. In December 1998, for instance, the government of Ontario, representing over a third of Canada's citizens, announced that it will issue identity certificates to its 11 million residents. Identrus [376], a joint venture set up in October 1998 by eight international banks, issues identity certificates for business-to-business electronic commerce. Equifax Secure, a division of Equifax (one of the three major U.S. consumer credit bureaus), in May 1999 announced an identity certificate service that matches information provided by individuals against data from Equifax Credit Information Services and other consumer and business information sources, to establish identity in real time. Strassman and Atkinson [363] propose that the U.S. Department of Motor Vehicles issue identity certificates.

identity-based PKIs by Geer [182], Kaufman Winn [227], Gladman, Ellison, and Bohm [187], Guida [200], and Ellison and Schneier [148].⁷)

- As Garfinkel [181, Chapter 4] explains, the approach of creating a society in which every person can be held accountable for his or her own actions by replacing *anonymity* (i.e., the privacy of identity) with absolute identity is fundamentally flawed. Identity will have to be established on the basis of legacy paper-based systems, and thus will inherit their insecurity. Also, identities may erroneously or maliciously be swapped or forged. Criminals who manage to steal identity certificates or to assume the identities of unwitting people will be able to misuse certificates in cyberspace on a global scale, while their victims take the blame. Punishment of the wrong individuals will make others reluctant to participate.
- In PKIs in which communicating or transacting parties have not established a prior relation, certificate verifiers will primarily be interested in the privileges and other non-identity attributes of certificate holders. If an individual's certificate includes all the attributes that a verifier needs to know in order to locally decide what action to take, many of the drawbacks of central database look-up are overcome.

Placing the data that would otherwise be listed in central database entries into attribute certificates is most natural in closed PKIs. A closed PKI is a PKI which has one issuer and clear contractual relationships between the issuer, certificate applicants, and verifiers. Closed PKIs are much more viable than open PKIs, where each certificate serves to establish authenticity in a potentially unbounded number of applications, since it is much easier to determine the risks and liabilities in a closed PKI. Moreover, organizations typically are not willing to let others issue certificates on their behalf, for commercial and liability reasons.

An early proposal for attribute certificates is due to Brands [54], in 1993. This proposal aims to protect the privacy of certificate holders, and forms the basis for many of the techniques that will be developed in this book. Conceptually, it builds on paradigms developed by Chaum [87, 88, 93, 107] in the period 1985–1992. Chaum advocated the use of credentials, which he defined [93] as “statements concerning an individual that are issued by organizations, and are in general shown to other organizations.” Chaum's credentials are not attribute certificates, though; they are digitally signed random messages that do not include a public key. For a discussion of the drawbacks of Chaum's approach, see Section 1.2.2.

In 1996, Blaze, Feigenbaum, and Lacy [33] also argued in favor of attribute certificates that do not reveal identity. Their focus is not on digital certificates, though, but on the design of a trust management system (called PolicyMaker) that enables

⁷Be warned that several of the fears and doubts that Ellison and Schneier [148] raise are in no way specific to PKIs.

verifiers to make decisions when presented with attributes and a request for access to a service. (See Blaze, Feigenbaum, Ioannidis, and Keromytis [32] for details of PolicyMaker and KeyNote, a related trust management system designed specifically for making Boolean decisions based on attribute certificates.) A similar trust management system is REFEREE [116], which forms the basis of the DSig [115] initiative of the World Wide Web Consortium; DSig is a proposed standard format for making digitally-signed, machine-readable assertions about a particular information resource. All these developments are orthogonal to, and can be used in conjunction with, the techniques that we will develop in this book.

A standardization effort for attribute certificates is the Simple Public Key Infrastructure [151] (SPKI). SPKI rejects not only the identity focus of the X.509 framework, but also its use of global names and its hierarchic certification structure; see Ellison [147, 149] for details on the SPKI design philosophy. In April 1997, SPKI merged with the Simple Distributed Security Infrastructure [324] (SDSI). SDSI is a PKI proposal by Rivest and Lampson; it is based on local names spaces, and centers around public keys rather than individuals. SPKI/SDSI 2.0 [152] combines the SDSI local names spaces and the SPKI focus on attribute certificates.

Tokeneer [319, 320], a PKI proposal by the NSA, heavily relies on attribute certificates as well, mainly because in federal agency applications immediate connectivity to a trusted authentication server is not always possible.

In 1997, VeriSign announced that it would personalize its digital ID's with a zip code, age, gender, and personal preferences, to facilitate integration with the Open Profiling Standard [209] (OPS). OPS was announced in November 1997 by Netscape, Firefly Network, and VeriSign as a framework for the automated transport of personal data of individuals to Web sites. The idea of OPS is that an individual enters his or her personal data once, after which it is stored in the form of a Personal Profile in encrypted form on his or her personal computer. Some or all of the personal data in a Profile may be digitally certified. A set of rules is then used to determine how and when the data can be disclosed to online services. In 1998, the Platform for Privacy Preferences [393] (P3P) of the World Wide Web Consortium subsumed OPS. P3P allows Web sites and visitors to automatically negotiate a degree of privacy, based on the privacy practices of the Web site and privacy preferences specified by the individual in his or her browser.

Several PKI proposals use signed attribute objects that are in fact not true attribute certificates, because they do not bind attributes to a public key. This approach is followed in X.509v3 extensions, and has been adopted amongst others by Netscape's Transport Layer Security (TLS) 3.1 and X9.57 of the American Bankers Association. An X.509v3 "attribute certificate" has the same syntax as an X.509v3 certificate, but has a null public key; to prevent replay, it has an embedded link to a standard X.509 identity certificate, the public key of which is used for authentication. A key holder may have multiple of these signed attribute objects associated with the same identity certificate. Advantages of this approach are that the attributes do not increase the size

of the identity certificate, and attributes can be refreshed independently of the identity certificate. RSA's PKCS #6 [330] embeds X.509 certificates into a structure that adds additional attributes before the whole package is signed, to provide backward compatibility with X.509 certificates; the resulting structure is a genuine attribute certificate.

Note that the validity period of an attribute certificate may not exceed that of the attribute with the shortest validity period. For instance, if an attribute specifies the age of a person, then any certificate in which that attribute is specified should not have an expiry date that extends beyond the person's next birthday. In this particular example, the problem can be removed by encoding the date of birth instead of age, but this is not always possible. In other words, there is an incentive to use *short-lived* certificates (i.e., certificates with short validity periods).

1.1.5 Certificate revocation and validation

Certificates are valid until they expire, unless they are revoked beforehand. Many things can happen that require the revocation of a certificate. For example, the secret key may be lost or irreversibly destroyed, the certificate holder may cease operation, the certificate holder's identifier may need to be updated due to a name change, one of the (other) attributes in the certificate may have become invalid, or the secret key may have been compromised. It is not necessarily the certificate holder who desires to revoke a certificate. For example, when a company fires an employee, it is often necessary to revoke all his or her access privileges. Also, a certificate holder who uses a *limited-show certificate* (e.g., a discount coupon or a public transit ticket) more times than allowed must be stopped from continuing the fraud.

While revocation is an exceptional circumstance, the task of verifiers to check the revocation status of unexpired certificates unfortunately is not. They must either have the certificate status validated online (at the time of the communication or transaction) or regularly download a digitally signed update of a blacklist called the Certificate Revocation List (CRL). In both cases the status of certificates must be maintained by the CA (or by a special Revocation Authority). Note that the certificate revocation or validation data must itself be authenticated.

X.509v1 and PEM rely on the distribution of full CRLs. X.509v2 introduced the notion of delta-CRLs, which are in essence CRL updates. In X.509v3, the set of all issued certificates is subdivided into fragments that each have their own CRL; each X.509v3 certificate has a pointer to the CRL fragment that indicates its revocation status ("CRL Distribution Points"). See Perlman and Kaufman [300], van Oorschot, Ford, Hillier, and Otway [292], and Adams and Zuccherato [5] for related proposals.

As an alternative to the CRL approach of X.509v3, the PKIX working group is standardizing an online validation method, called the Online Certificate Status Protocol [270] (OCSP), for time-critical applications. ACES rejects the CRL approach altogether in favor of an online validation check, to facilitate a "pay as you go" busi-

ness model; the idea is that federal agencies pay a certificate validation fee each time they rely on certificates issued by commercial CAs for authentication.

Online certificate validation avoids the need for verifiers to manage their own versions of a CRL and to deal with certificates they are not interested in, but suffers from all the problems of the central database paradigm. One of the primary problems is scalability to large communities, not in the least because responses to queries must be authenticated by the trusted central database. In fact, in many PKIs (especially those with just one CA) it makes little sense to use digital certificates in combination with online certificate validation; organizations or the CA might as well keep copies of public keys on file. Improvements of the basic mechanism for online certificate validation have been proposed (see Kocher [235], Micali [268], Aiello, Lodha, and Ostrovsky [7], and Naor and Nissim [273]), but these do not remove the main problems.

Distribution of CRLs (or their updates), is more attractive in many respects, but creates a lag between the time a certificate becomes invalid and when it appears on the next CRL update. If validity periods are long, CRLs will grow and additional computing resources are needed for searching and storing them.

Either way, certificate revocation seriously reduces the finality of secure communications and transactions. In 1994, the MITRE Corporation [27] estimated that the yearly running expenses of an authentication infrastructure derive almost entirely from administrating revocation. Its cost estimates are based on the distribution of full CRLs, but would be similar for CRL updates, and probably worse for online certificate validation. Another serious problem is that revocation requires secure time stamping, because otherwise one can simply backdate signatures. (Likewise, verifiers need a secure clock to verify the expiry dates or validity windows of certificates.)

PGP leaves it to key holders themselves to notify their correspondents in case their keys are to be revoked, and relies on the revocation information to propagate. This approach works well in small communities, but is not workable in large-scale PKIs where key holders have transient relations.

In the currently prevailing certificate paradigm, certificates are *long-lived*; validity periods are typically in the order of many months or even years. The use of long-lived certificates is often taken for granted, presumably for historical reasons: getting multiple copies of a paper-based certificate whenever desired is not a feasible option. However, with today's computers and electronic networks, getting 100 certificates is hardly less efficient than getting a single one, and one can always reconnect with the issuer to download a new batch of certificates. In many cases, issuing certificates with short validity periods is sufficient to deal with the revocation problem, as Kaufman [226] and others observed. Elaborating on this observation, Stubblebine [365] proposed to include recency information (validity windows) and freshness policies within certificates, to reduce the importance of timely revocation information. (See McDaniel and Jamin [260] for a related proposal.) Rivest [323] proposed to abolish certificate revocation altogether, by having the certificate holder

supply all the evidence needed by the verifier to check the validity of a certificate; freshness is achieved by showing a more recently issued certificate.⁸ SDSI 2.0 for this reason uses no revocation mechanism at all. Diversinet, which provides what it calls digital permits (similar to the X.509v3 construct for “attribute” certificates), also follows this approach.⁹ Clearly, the model of short-lived certificates fits well with many types of limited-show certificates. In Chapters 5 and 6 we will see that the paradigm of short-lived limited-show certificates has many other benefits over that of long-lived unlimited-show certificates.

Note that revocation is not needed when a secret key has been destroyed or its holder voluntarily ceases operation, assuming the certificate served only for authentication purposes and certificate holders act only on their own behalf.¹⁰

1.1.6 Smartcard integration

Everything discussed thus far applies to certificates implemented in *software-only computing devices*. These devices may be obtained on the open market and may be modified freely by their holders; they do not contain any tamper-resistant components that serve to protect the security interests of their issuer. Examples are desktop computers, notebooks, palmtop computers, cellular telephones, and smart watches. Software-only implementations have many advantages: mass-scale software is cheap (software needs to be written only once), can be manufactured in-house by any software producer, and is easy to distribute over the Internet and other networks (no need for physical transportation); any individual can issue his or her own certificates (low start-up cost); and, it is relatively easy to verify claims about the operation of the software.

Nevertheless, software-only implementations are not preferable in most PKIs. A distinct problem is theft. If the secret key of a certificate is generated and stored on a personal computer or the like, it is virtually impossible to prevent its compromise, loss, disclosure, modification, or unauthorized use. Gutmann [204], for instance, in 1997 found that “no Microsoft Internet product is capable of protecting a user’s keys from hostile attack,” due to design and implementation flaws. Shamir and van Someren [348] note that software run by an attacker (in the form of a virus or a Trojan horse, or simply from a floppy during lunch-time) may be able to rapidly detect cryptographic keys stored on a PC by scanning for data sections with unusually high entropy. Encrypting the secret key does not overcome the problem: encrypting it using a password is vulnerable to a brute force attack, and in any case the secret key must at some stage be in the clear to be usable.

⁸McDaniel and Rubin [261] argue that this approach is not preferable over CRLs in all circumstances.

⁹A Diversinet certificate attests to the binding between a public key and an anonymous identifier that can be uniquely linked to centrally maintained identity information and other attributes; see Brown [64].

¹⁰In contrast, a public key used for encryption should be revoked when the secret key is lost, but there is no satisfactory way to do this securely. Another secret key must be established to enable authenticated communication with the CA; this shifts the problem but does not overcome it.

Also, in many PKIs participants are not allowed to lend, share, or give away their certificates. Software-only devices cannot protect against this. Examples of *personal certificates* are driver's licenses, diplomas, subscriptions, electronic passports, and employee badges. To limit transferability, Lessig and Resnick [249] suggest to include location data (such as an IP address) in certificates or to make certificates traceable (so that abusers can be punished); both methods offer inadequate security, though, and destroy privacy. For some odd reason, the lending problem is not widely acknowledged, as witnessed by the limited amount of work done to protect against it.

A convenient way to protect against these and other risks is to store secret keys on a *smartcard*. This is a chipcard that contains a microprocessor that is capable of making arithmetic decisions.¹¹ Smartcards can process data in intelligent manners, by taking actions based on secret data that never needs to leave the card. Memory access and input/output are guarded against unauthorized access, and the card can disable itself after a false PIN has been entered several times. (Alternatively, the card can store an electronic representation of its holder's fingerprint, and match this against a fingerprint entered on a trusted card reader or directly on the card.) Tampering with a smartcard in order to get to its contents can set off an alarm in the card that blocks it or overwrites the memory contents with all zeros (a process known as zeroization).

Implementations based on tamper-resistant smartcards offer a multitude of benefits, many of which have systematically been overlooked by PKI researchers and developers:

- It is easy to protect smartcards against viruses and Trojan horses aimed at capturing their secret keys.
- If the smartcard has an access control mechanism (PIN, password, or otherwise), certificates cannot be shown by smartcard thieves and other parties not authorized by the card's legitimate holder.
- If the smartcard has an access control mechanism that scans a biometric of the card holder (e.g., his or her fingerprint), certificate holders can be prevented from lending their personal certificates.
- The tamper-resistance can prevent certificate holders from making copies of limited-show certificates. In software-only implementations, the only way to protect against fraud with limited-show certificates is to require online clearing with a central party for all transactions.

¹¹The term smartcard is used differently by different organizations. The definition of smartcards used by ISO 7816-1 of the International Organization for Standardization, and applied by the Smart Card Forum, includes memory cards. It stands to reason that organizations aiming to promote and commercialize a new technology prefer to use a single term as broadly as possible. However, a memory card can hardly be said to be "smart," even if it has hardwired logic.

- If the smartcard has a tamper-resistant internal clock, the burden of checking the expiry date of a certificate can be moved from the certificate verifier to the smartcard, with the latter refusing to (help) show a certificate if the date is outside of its validity period. This avoids the need for certificate verifiers to run a secure clock that cannot be reset by an attacker to a time within the validity period.
- If the smartcard has a tamper-resistant internal clock, it can add a timestamp to any digital signatures made when showing certificates.
- If the smartcard has a tamper-resistant internal clock, it can limit the number of certificate showings within a given time period; this may be desirable for certain types of certificates or CA liability arrangements.
- Certificates that specify negative qualifications can be discarded only by muzzling the smartcard. The smartcard can then enter into suspension mode, so that its holder can no longer show any certificates.
- More generally, the smartcard can locally decide whether its holder may engage in certain transactions, and can prevent undesired behavior.
- Certificate holders cannot fall victim to extortioners in cyberspace who (possibly anonymously) extort them into transmitting their certified key pairs to them; an extortioner cannot reasonably expect his or her victim to be able to break the tamper-resistance of the smartcard.
- The smartcard can prevent its holder from helping remote parties (over a radio link or the like) to gain access to services for which they do not have the proper certificates themselves. (Details will be provided in Section 6.5.3.)
- The smartcard can do internal book-keeping in the interest of its issuer, such as keep track of an electronic cash balance. It could even keep a log of all transactions, which could be inspected by law enforcement agents that have a court order. Of course, any reliance on smartcard book-keeping is acceptable only when the damage that can result from tampering with a card's contents is outweighed by the cost of breaking the card's tamper-resistance.
- The vulnerability of secret keys stored in software-only devices makes it difficult to reliably associate a digital signature with a particular individual. This makes the legal status of digital signatures doubtful, which in turn hampers the progress of electronic commerce. Tamper-resistant devices for certificate holders help give digital signatures a firm legal grounding.
- The need to rely on revocation mechanisms greatly reduces. The latency and scalability problems of software-only CRL distribution are largely overcome:

the risk of theft is minimized (assuming the presence of an access control mechanism), it takes significant time to physically extract keys from stolen smartcards, and a card holder's capabilities can be revoked by taking back the smartcard. Consequently, there should rarely be a need for online certificate validation and the expected size of CRLs (assuming certificates are used only for authentication) is zero. In many closed smartcard-based PKIs, the validity of certificates can even be made the liability of the CA.

- Smartcards offer portability.
- Smartcards have a clear psychological advantage over software-only implementations.

In sum, there are many reasons to prefer smartcard-based implementations over software-only implementations. Of course, any other tamper-resistant hardware devices may be used instead. Although smartcards are a natural choice in many situations, other embodiments may be more appropriate for securely implementing an internal battery, a clock, a live access control mechanism, or other desirable features. One interesting alternative is the iButton¹² produced by Dallas Semiconductor. For concreteness, however, throughout this book we will always refer to smartcards whenever there is a need for tamper-resistant devices for certificate holders.

One of the most successful smartcard-based PKI implementations is Fortezza. Other smartcard-based PKIs relying on identity certificates are Chip-Secure Electronic Transactions [20] (C-SET) and version 2.0 of SET. Pretty Good Privacy and Schlumberger Electronic Transactions in April 1997 announced a strategic alliance for the development and marketing of PGP-enhanced smartcards. VeriSign in November 1999 announced to bundle its Class 1 Digital IDs with Litronic's NetSign, a technology that integrates Netscape Communicator with smartcards for increased portability. Identrus and GTE CyberTrust have announced similar plans. In June 1999, Gemplus Software (a division Gemplus) announced GemSAFE Enterprise, a smartcard solution intended to integrate seamlessly with X.509 and other popular certificate standards and products. In the same month, eFed, Entrust Technologies, and NDS Americas demonstrated a smartcard-based PKI procurement solution for the federal marketplace, based on X.509 certificates. More recently, in November 1999, the U.S. Department of Defense announced that it will use smartcards in its PKI plans, in order to securely identify military, civilian and contractor employees when they gain access to its buildings, computers, Internet and private networks, and so on; by the end of 2000 a roll-out of some 4 million smartcards will begin.

¹²The iButton is a 16mm computer chip housed in a stainless steel can, which can be affixed to a ring, key fob, wallet, watch, or badge. The cryptographic iButton version contains a microprocessor and high-speed 1024-bit coprocessor for public-key cryptographic operations, runs a Java virtual machine, has 6 kilobytes of SRAM that will zeroize its contents in case of an attempted physical compromise, and transfers data at up to 142 kilobits per second.

ACES also foresees an important role for identity certificates stored on smartcards. According to the U.S. Federal Card Services Task Force [161], the goal of the U.S. government is “to adopt an interoperable multi-application smart card that will support a wide range of governmentwide and agency-specific services. This goal sets the target for every federal employee to carry a smart card that can be used for multiple purposes – travel, small purchases, identification, network and building access, and other financial and administrative purposes– by the year 2001. [...] This plan calls for a smart card based extended ID authentication function to support multiple applications based on public key technology using the standard X.509v.3 digital certificate and authentication framework as an operating model.”

One of the few smartcard-based PKI proposals specifically designed for attribute certificates is NSA’s Tokeneer [320]. As its authors [319] note, there is an urgent need to minimize the amount of data transferred both to and from the smartcard:

“In the case where the certificate is stored in a token (such as a smartcard), storage area may be a critical factor. Each certificate requires at least one signature. The size of the signature depends on the exact signature algorithm being used (128 bytes in the case of a 1024 bit RSA signature). Adding other data fields and ASN.1 encoding overhead, each certificate can be on the order of several hundred bytes. [...] In systems where I/O throughput is a factor (especially in smartcard based systems where I/O may be limited to 9600 baud/second) the data size of the certificate may be a concern. Separating the certificate into several certificates will be beneficial only if transfer is limited to one certificate per service request. Separating attributes into several different certificates may also be detrimental to overall system performance if multiple certificates are required.”

The computation, communication, and storage burden is frequently cited as one of the main reasons why smartcard implementations of certificate mechanisms are stalling. Indeed, current proposals all require sophisticated smartcards with large EEPROM¹³ and cryptographic coprocessors. Efficiency considerations are of prime importance in smartcard implementations, especially for short-lived and limited-show certificates. The addition of complex circuitry and software is expensive and can easily lead to new weaknesses in the internal defense mechanisms. Also, with smartcard components already cramming for space, adding circuitry adversely affects reliability. A “dead” card inconveniences and frustrates its holder, and can have dramatic consequences in medical and other applications.

The techniques that we will develop in this book overcome all these problems.

¹³EEPROM, or Electrically Erasable Programmable Read Only Memory, is non-volatile memory that enables data to be erased by discharging gates electrically.

1.2 Privacy issues

The efficiency and security problems of certificate revocation and smartcards are not the only shortcomings of the current proposals for PKIs and certification mechanisms. In this section we examine the privacy dangers, and discuss the meager efforts that have been undertaken to protect privacy. On the basis of this analysis we then list desirable privacy properties for digital certificates and PKIs.

1.2.1 Privacy dangers

As defined by Westin [387], (information) *privacy* is the claim of individuals, groups, or institutions to determine for themselves when, how, and to what extent information about them is communicated to others.¹⁴ It is a fundamental postulate of this book that if the current visions about the global PKI (i.e., the collection of all regional, national, and international PKIs) turn into reality, then everyone will be forced to transact and communicate in what will be the most pervasive electronic surveillance tool ever built. (*Surveillance* is the act of systematically monitoring, tracking, or assessing the actions of individuals.)

To apprehend the magnitude of the privacy problem, consider the following aspects:

- All the communications and transactions of an individual in a PKI can be linked on the basis of his or her identity certificates. In this manner, dossiers can automatically be compiled for each individual about his or her habits, behavior, movements, preferences, characteristics, and so on. Many parties enjoy this dossier forming capability:
 - The CA sees the certificates it issues, and typically sees them again once they are shown or at a later moment. This enables the CA to trace and link all communications and transactions of each key holder. Reasons why the CA may get to see the certificates that are shown to verifiers include: the verifiers in the PKI may belong to the same entity as the CA (closed system); verifiers may be incited to deposit a copy of the transcript of each transaction they engage in (e.g., to enable detection of certificate forgery or fraud with limited-show certificates, or to support commercial goals); and, verifiers may resort by default to online certificate validation with the CA.¹⁵

¹⁴Westin's definition is frequently cited in the academic literature and in court decisions, and forms the basis for the U.S. Privacy Act of 1974 (for an overview, see, e.g., the Office of Technology [285, Chapter 1]) and similar legislation in many other developed countries around the world (see the Global Internet Liberty Campaign [153] and Rotenberg [329]).

¹⁵Today's PKIs seldomly require certificate verifiers to deposit their transcripts to the CA, but this can be expected to change as the awareness of the security benefits grows.

- Each verifier can store all the certificates that are presented to it, and can link them on the basis of their key holder identifiers, public keys, or CA signatures. Different verifiers can exchange and link their data on the same basis. Developments are underway to streamline the latter process. For instance, the Information Content & Exchange protocol (ICE), announced by the ICE working group in June 1998, provides organizations with a standardized automated method for exchanging personal data obtained through P3P and other mechanisms. Another planned standard, the Customer Profile Exchange (CPEX), announced in November 1999, is intended to take automated exchange and integration of personal data to an even further level. The eXtensible Markup Language (XML), designed by the World Wide Web Consortium, will play an important role in these initiatives.
- In case the communications or transactions of key holders are not securely encrypted, wiretappers see the same information as verifiers. If they can wiretap the certificate issuing process as well, they can learn everything the CA knows. One particularly nasty aspect is that the certified public key of at least one of the two parties in a transaction or communication is always sent in the clear to bootstrap a secure session.

Parties that actively monitor the Internet and other telecommunication infrastructures, or at least have the capability to do so, include government agencies (international wiretapping efforts include Echelon¹⁶ and Enfpol¹⁷), non-profit organizations (such as the Internet Archive, which stores over 14 terrabytes of information gathered from news groups, Web pages, and other publicly accessible Internet sites), and commercial enterprises (e.g., Internet routers, Internet service providers, and the commercial offshoot of the Internet Archive). The U.S. Communications Assistance for Law Enforcement Act [145, 159, 160, 286] and similar legislation [359] in other countries require the telecommunications industry to build wiretapping capabilities into their infrastructures.

- The CA can trivially link each dossier to the identity of each individual. For verifiers and wiretappers, linking dossiers to identities typically is straightforward as well: either separate entries or the aggregated contents of a dossier reveal the identity, or the match can be made in another way (e.g., on the basis

¹⁶Echelon is an international surveillance system that taps into most of the world's non-military satellite, radio, and land-based communications systems. It is operated by the United States in cooperation with Great Britain, Canada, New Zealand, and Australia. Echelon systematically scans e-mail, fax, cellular, telex, and telephone communications for keywords, to identify and extract messages deemed of interest. See [75, 206, 207, 308] for details.

¹⁷The Council of the European Union and the FBI have been cooperating since 1992 on a plan for intercepting all mobile phone calls, Internet communications, and fax and pager messages in Europe. See [129, 262, 344, 362] for details.

of voice or facial recognition or by tracing the source of an Internet connection). Sending your digital certificate offers no more privacy than sending your Social Security number or credit card number.

Worse, all the dossiers that are compiled by linking and tracing the actions of participants in one PKI can be tied to the dossiers compiled in other PKIs. In the original X.509 proposal, each key holder would be assigned a globally unique identifier, providing a highly convenient method to link the actions of key holders across different PKIs. The SPKI authors [151] rightfully note that this “would constitute a massive privacy violation and would probably be rejected as politically impossible.” The use of local names, as in SDSI, makes it more difficult to link transactions across different PKIs, but clearly with today’s network resources and linking power the barrier that is raised is low; after all, local names and other kinds of identifiers are strongly correlated to true names. In any case, different local names of the same individual can all be linked when CAs cooperate.

- Attribute certificates worsen the problem, since the dossiers that CAs, verifiers, and wiretappers can compile are often even more intrusive. Attribute certificates that do not specify explicit identifiers can be linked and traced as easily as identity certificates, on the basis of their public key or the signature of the CA. (The same holds for X.509v3 “attribute” certificates and Diversinet’s digital permits, in spite of the latter’s claim¹⁸ that its digital certificates assure “total anonymity and privacy by separating authorization credentials into permits;” protection against unsophisticated wiretappers is a far cry from privacy, and can simply be achieved by line encryption.) Also, each CA gets to learn all the attributes of each certificate applicant, because otherwise it cannot or will not issue a certificate. Some CA service providers, such as Thawte Certification, are promoting PKI models whereby a single CA validates all the attributes of a certificate applicant, to avoid cumbersome verification procedures; this further increases the power of the CA.

Ellison [150] states: “Because SPKI certificates will carry information that, taken together over all certificates, might constitute a dossier and therefore a privacy violation, each SPKI certificate should carry the minimum information necessary to get a job done.” Indeed, SPKI certificates are not programmable; they have 5 exactly fields (Issuer, Subject, Delegation, Authorization, Validity Dates). A noble sacrifice, but one that one would prefer to avoid; in many PKIs, everyone benefits when more attributes can be encoded.

- Any digital signatures that are made by certificate holders can be added to their dossiers. They form self-signed statements that cannot be repudiated,

¹⁸At www.dvnet.com/about_us/what_we_do.htm, last checked March 30, 2000.

proving to the whole world who is the originator of a message. As Directorate-General XIII of the European Commission [137] notes, “Digital signatures could even bring significant law enforcement benefits as they allow for example messages to be attributed to a particular reader and/or sender.” In the words of Walsh [386], a former deputy director-general of the Australian equivalent of the NSA, “If you ever allow people to get near authentication keys you’ll corrupt the administration of justice.”

In a similar manner, anyone who gets to see a digital certificate, either by wiretapping a communication or by consulting online certificate repositories or CRLs, has convincing evidence that the identity and any other attributes signed by the CA belong together. Obtaining this information is often perfectly legitimate even for outsiders; many schemes (e.g., X.509 and PGP) store certificates in mail servers or other public depositories. The American Bar Association [9] states: “Publication of a certificate which has not been accepted by the subscriber may disclose an identification, business relationship, or other fact which the purported subscriber wishes to keep confidential, and may have a right to keep confidential under applicable privacy law.” Clearly, accepting subscribers have similar concerns.

- Any uniquely identifying data in a certificate (such as a key holder identifier, the public key, or the CA’s signature) can be misused to deny a key holder access to PKI services, and to block his or her communication attempts in real time. For example, blacklists can be built into Internet routers. Similarly, transaction-generated data conducted with target public keys can be filtered out by surveillance tools, and electronically delivered to third parties for examination or immediate action. More generally, entire groups can be discriminated against on the basis of attributes encoded into their certificates.
- Revocation mechanisms cause additional privacy problems. CRLs (or their updates) are distributed to all verifiers, and potentially to anyone who requests them. In this manner, entities can collect data about key holders they have never communicated or transacted with. Furthermore, the CA can falsely add public keys to its CRL, to block the communications and transactions of targeted certificate holders. (The methods of Kocher [235], Micali [268], Aiello, Lodha, and Ostrovsky [7], and Naor and Nissim [273] do not protect against false claims of the CA that a certificate has been revoked; they protect merely against a misbehaving Revocation Authority that gets authenticated revocation information from an honest CA.¹⁹)

¹⁹An improvement would be for blacklists to list the (serial numbers of) certified public keys together with a “suicide note,” a revocation message signed using the secret key. This approach, which is applied in PGP, cannot be used when the secret key is lost, while preparing a suicide note in advance is not an adequate solution.

Online certificate validation services are even worse: they allow anyone to verify not only negative data but also positive data, and enable the Revocation Authority to falsely deny access to certificate holders and to learn in real time who communicates with whom. This cripples the privacy of certificate verifiers as well.

- The integration of smartcards exacerbates the privacy problem. Moreno, the inventor of the first generation of smartcards in the early seventies, warned that smartcards have the potential to become “Big Brother’s little helper.” The tamper-resistance of a smartcard shields its internal operations from its holder. It is difficult or even impossible to verify that a card does not leak personal data and other attributes that may be stored inside the card. Leakage may take place by exploiting the van Eck effect,²⁰ by sending out or receiving radio signals, by sending along additional data when engaging in a protocol, by encoding information in message fields or random numbers specified in the protocol, by timing the delay before transmitting a message, or by halting at a specific step of a protocol. Also, the smartcard issuer (typically the CA from which the holder obtains certificates) can program the card in such a manner that it will lock its holder out of services upon receiving a signal from (the terminal of) a certificate verifier or another party. See Section 6.1.1 for details.

As more and more personal data is stored inside smartcards, individuals will be misled into believing that they have control over their own data. Consumer protection agencies such as the Privacy Commissioner of Canada [287] and the Privacy Commissioner of the Commonwealth of Australia [311] have already expressed great concern. See also Cavoukian, Johnston, and Duncan [79], Clarke [117], Connolly [120], Fancher [158], Schwartz [343], and Wright [394].

Since the surveillance of automated transaction systems is more surreptitious than wiretapping, it has even greater potential to law enforcement and intelligence agencies. Transaction-generated data trails can readily be picked up by computers, stored in databases, searched for patterns of activity, processed to distill profiles, and merged and matched with census data, credit report data, postal codes, car registrations, birth certificates, and so on. Moreover, transactions need not be monitored in real time; once stored, data trails are permanent for the record and can be examined at any time. Indeed, as NIST’s FPKI chairman Burr [68] notes: “Archives provide a long term storage of CA files and records. The life time of CAs may be relatively short. But it may be important to verify the validity of signatures on very old documents. CAs must make provisions to store the information needed to verify the signatures of

²⁰Microprocessors, keyboards, computer monitors, serial cables, printers, and other peripheral devices all emit electromagnetic radiation that passes over large distances and through solid walls, and that can be remotely captured and viewed; see van Eck [379] for the (purposely incomplete) paper that brought this phenomenon to the attention of the public.

its users, in archives that will be able to make the data available at a much later date, perhaps several decades later.”

Hardly surprising, individuals are feeling increasingly threatened. A survey conducted in April 1998 by Louis Harris & Associates and Westin for Privacy & American Business and Price Waterhouse found that 87% of American computer users are concerned about threats to their personal privacy. The threats do not merely pertain to abuse by the private sector. Indeed, as Singleton [356] points out, “Although both private and government databases can be abused, the abuse of government databases poses a more serious threat for one reason: government controls the courts, the police, and the army.” Commercial organizations have a commercial self-interest in protecting the privacy of individuals, and often are less interested in the behavior of identified persons than government agencies. Moreover, the surveillance technologies used by governments are typically more advanced and covert than those used by the private sector. In an influential study conducted in 1996, the U.S. National Research Council [278] “acknowledges the concerns of many law-abiding individuals about government surveillance. It believes that such concerns and the questions they raise about individual rights and government responsibilities must be taken seriously. It would be inappropriate to dismiss such individuals as paranoid or overly suspicious.”

ACES and other government PKIs without privacy-protection measures intrude even more on privacy than the national ID cards that many developed countries are considering in order to combat tax evasion, social security fraud, illegal immigration, insurance fraud, fraudulent work authorization documents, and so on. These plans have already led to public outcry in the United States, Great Britain, Canada, New Zealand, Australia, and other countries. See Privacy International [312] for an overview of (proposed) national ID cards around the world.

1.2.2 Previous privacy-protection efforts and their shortcomings

Surprisingly, the issue of privacy in PKIs has received virtually no attention. Most certification technologies and standardization efforts do not deal with the issue at all, or only allow users to encrypt their communications and transactions at the transport layer. Confidentiality is a weak privacy measure, though. As Baker [18], then chief counsel for the NSA, remarked: “The biggest threats to our privacy in a digital world come not from what we keep secret but from what we reveal willingly. [...] Restricting these invasions of privacy is a challenge, but it isn’t a job for encryption. Encryption can’t protect you from the misuse of data you surrendered willingly.”

The European Commission [137] recommends that individuals be allowed to obtain digital certificates that specify a pseudonym, unless the law specifies that true names must be used. This approach is supported by the OECD Cryptographic Policy Guidelines [295] and by the European Privacy Directive [157].²¹ *Pseudonymous cer-*

²¹The European Privacy Directive is an extensive set of privacy guidelines established in 1995 by the

tificates are recommended also by, amongst others, Birch [31], Gladman, Ellison, and Bohm [187], Hill and Hosein [211], and Standards Australia [361]. PKI efforts that provide for pseudonymous certificates include X.509v3 (certain agents may protect their identity through the use of role-titles; “residential persons” do not enjoy this capability, though), PEM (pseudonymous certificates can be retrieved through so-called Persona CAs), SDSI (its free-form identity certificate syntax allows the specification of pseudonyms), VeriSign’s Digital IDs (VeriSign issues identity certificates with anonymous identifiers), and PGP (users can specify false email addresses). Clearly, pseudonymous certificates that can only be obtained by certificate applicants who identify themselves to the CA offer no better privacy than Social Security numbers and credit card numbers; at least their issuer can readily follow them around and trace them to the identity of their holder. The alternative of anonymous registration of certificate applicants offers better privacy, but unfortunately suffers from serious drawbacks:

- It may be very difficult to register without being identified by cameras or personnel, or via one’s IP address. Typically it is easier to realize an anonymous channel when showing a certificate than when retrieving the certificate, especially in the physical world. Identification in any one interaction with the CA results in one’s communications and actions becoming traceable.
- Anonymous registration does nothing to prevent linkability of all the communications and transactions of certificate holders. To overcome this problem, each certificate must be retrieved anonymously on a separate occasion, with a significant random delay between each retrieval to prevent linkability by the CA. Not only is this impractical, it also prevents certificate applicants from building a reputation with the CA. In particular, the CA cannot distinguish frequent from infrequent certificate applicants and cannot lock out fraudulent users.
- Many types of certificates, in particular personal certificates, are issued only to identified applicants. Even if the CA certifies only personal attributes that do not identify the certificate holder, such as age and marital status, often the only way for the CA to verify the attributes is by establishing the applicant’s identity and using this to look up the attributes in a trusted database.
- Non-repudiation and recovery of lost or stolen certificates are hard to implement.
- It is hard or even impossible to protect against basic forms of fraud, including unauthorized lending, copying, and discarding of certificates. Many applica-

European Union, requiring the 15 member states to harmonize their national laws to protect personal information. It has taken effect in October 1998, and applies to commercial as well as governmental information processing. Numerous non-EU countries have adopted privacy policies that are in alignment with the European approach; see Davies [130].

tions require mechanisms through which misbehaving key holders can be identified, so that they can be locked out from further participation and possibly be sued for damages. Traceability also serves to discourage those contemplating fraud. Anonymously issued certificates do not offer these security measures, in fact the CA cannot even contain the damages due to fraud.

- It is impossible to purchase certificates from the CA while remaining anonymous, unless one can pay using hard cash or a privacy-protecting electronic cash system.

These drawbacks make anonymous registration of certificate applicants a highly undesirable course of action in the vast majority of PKIs (unless one resorts to the cryptographic techniques that will be developed Section 5.5.1).

OPS and P3P are erroneously hailed as technical solutions to the privacy problem. They make it much easier to obtain personal data from individuals, and can be abused by service providers to turn away or discriminate against persons who do not want to disclose their identity and other personal data; in this manner they compel people to give up their privacy. The Data Protection Working Party of the European Union [214] criticized P3P on the grounds that it seeks to “formalise lower common standards,” and that it “could mislead EU-based operators into believing that they can be discharged of certain of their legal obligations (e.g. granting individual users a right of access to their data) if the individual user consents to this as part of the on-line negotiation.”

VeriSign issues certificates that contain encrypted attributes that can be unlocked only by verifiers that meet the qualifications necessary to receive the required decryption key from a trusted third party. This encryption measure does not reduce the surveillance capabilities of the most powerful parties in any way, nor does it prevent anyone from linking and tracing all the actions of each certificate holder.

Another attempt to protect privacy is for the CA to digitally sign (salted) one-way hashes of attributes, instead of (the concatenation of) the attributes themselves. When transacting or communicating with a verifier, the certificate holder can selectively disclose only those attributes needed.²² This generalizes the dual signature technique applied in SET [257]. Although certificate holders now have some control over which attributes they reveal to verifiers, they are forced to leave behind digital signatures. Furthermore, they are seriously restricted in the properties they can demonstrate about their attributes; Boolean formulae, for instance, are out of the question. Worse, nothing prevents the CA and others from tracing and linking all the communications and transactions of each certificate holder.

Ellison [150] states: “Because one use of SPKI certificates is in secret balloting and similar applications, an SPKI certificate must be able to assign an attribute to a blinded signature key.” *Blind signatures* are a concept introduced by Chaum for the

²²Lamport [244] proposed this hashing construct in the context of one-time signatures. When there are many attributes, they can be organized in a hash tree to improve efficiency, following Merkle [267].

purpose of anonymous electronic cash [91, 92, 96] and credential mechanisms [87, 93, 107].²³ While they can be used to overcome several of the drawbacks associated with anonymous registration of certificate applicants, they are not suitable to design attribute certificates:

- Because users can fully blind all the certificate data they obtain from the CA, it is not possible for the CA to encode person identifiers that must be disclosed in certain circumstances but may remain hidden in others.
- For the same reason, the CA cannot encode expiry dates into certificates. At best it can use a different signing key for different certificate issuances, and declare in advance when each issuance will become invalid.
- More generally, the blinding prevents the CA from encoding attributes into a certificate. It is possible to represent each combination of attribute values by a different signing key of the CA, but this seriously limits the number of attributes that can be encoded and their value ranges, and moreover an exhaustive list of the meanings associated with all possible signature types must be published in advance.²⁴
- Certificate holders cannot selectively disclose their attributes, since verifiers need to know which public key of the CA to apply to verify a certificate. A blind signature guarantees absolute anonymity and untraceability; it does not enable one to negotiate a degree of privacy.
- Chaum's credentials must all be created by the same party; different organizations who wish to issue credentials must all rely on this central party to do the factual issuance.
- Chaum's cash and credential mechanisms rely heavily on a real-time connection with a central party during each transaction. The requirement of online clearing and central database look-up during each transaction strikes against the philosophy behind digital certificates.

²³A blind signature scheme enables a receiver to obtain a signed message from a signature issuer in such a manner that the signed message is statistically independent from the issuer's view in the protocol execution.

²⁴In Chaum's proposal for RSA signatures, the signer uses the same RSA modulus but a different public exponent v_i (specifically, the i -th odd prime in sequence) to issue blind signatures that represent the i -th message in a public list. A variation would be for the CA to declare merely that the signing of a particular message requires as the public signature exponent the nearest prime exceeding the message (when viewing its binary representation as an integer), but this is impractical as well: a coding scheme must be applied to ensure that messages have sufficiently large Hamming distance, and generating and verifying a (new) signature type in addition to the normal workload requires (possibly a great many) applications of a primality test.

- Chaum's credentials do not have a built-in secret key to authenticate actions performed with them, and so the problem of replay arises. To prevent replay, the credential holder must authenticate the showing of a credential to an organization by applying the secret key of a digital pseudonym established (previously or at the same time) with that organization.
- Certificates that contain unfavorable attributes (i.e., attributes that the holder would prefer to hide, such as a mark for drunk driving) can simply be discarded by their holder. In many cases it is impractical to require all certificate holders to obtain and show attributes that indicate the absence of unfavorable attributes.
- Blind signatures cannot prevent or discourage the unauthorized lending of certificates. This drawback by itself renders the blinding technique useless for the majority of certificate applications.
- An extortioner can digitally extort certificates by forcing the victim to retrieve certificates for which the extortioner instead of the victim supplies the blinded messages. The victim merely signs the certificate request and passes the responses of the issuer on to the extortioner. The extortioner can subsequently at his or her leisure show the certificate. At no stage is there a need for physical proximity or a physical communication or delivery channel,²⁵ and so the extortioner can remain untraceable throughout.
- Chaum's smartcard techniques require smartcards with cryptographic coprocessors and plenty of memory, to guarantee that the required operations can be performed in reasonable time.
- More seriously, Chaum's smartcard techniques are not secure. Since attributes are encoded by the smartcard rather than by the CA, physical compromise of a smartcard enables its holder to forge attributes, and to lend, give away, or distribute copies of new certificates. The CA cannot trace fraud, and containment can only be accomplished by suspending the entire system. See Section 6.2.2 for details.
- Chaum's technique [90, 98] for one-show certificates makes high computation and communication demands on both the issuer and the receiver. Furthermore, it does not extend to limited-show certificates, does not admit zero-knowledge proofs, and cannot be migrated to a setting with smartcards without further

²⁵The extortioner can transmit from behind a computer that is part of a network located behind a firewall, use some a computer at an Internet cafe or a public library, deploy anonymous remailers, pseudonymous remailers, or Mixmaster remailers (see Goldberg, Wagner, and Brewer [190] for an overview), or use anonymity or pseudonymous services like Janus [35], Babel [203], Crowds [322], Freedom [189], and Onion Routing [368]. Also, most Internet access providers dynamically assign IP addresses to each client session; these identify only the host name of the computer used by the access provider to establish the session.

degrading efficiency. Also, it does not give certificate holders the ability to selectively disclose information about attributes.

In spite of their unsuitability for digital certificates and PKIs, Chaum's paradigms and methodologies provide valuable insight in how the privacy problems of PKIs may be overcome.

The certification mechanisms that will be presented in this book overcome all the problems mentioned in this section.

1.2.3 Desirable privacy properties

In many situations there is no need to disclose one's identity. For example, when a police agent stops an individual for speeding, all that the officer normally needs to know is whether the individual has a valid driver's license. When requesting entry to a gaming parlor it suffices to demonstrate one's age or year of birth. Likewise, a county database service may merely need to know that someone requesting a file is a resident. More generally, in many PKIs the use of identity is no more than a way of adding a level of indirection to the verifier's authentication algorithm; recall Section 1.1.3.

Even in PKIs where one would expect only identified actions, there may be a need for the ability to hide identity; the MITRE Corporation [27, page D-14], for instance, points out an application where an undercover FBI agent must file a report from a remote computer.

In today's computerized world we cannot expect others to protect our privacy. In order to protect privacy, we must operate under the following assumptions:

- (Persistence) Whenever data about certificate holders can be collected, it will be collected and stored indefinitely (if only because not collecting data that can so easily be collected must be considered a waste of resources). Every piece of information that is electronically submitted is there for the public record, even though the sender rarely intends the data to endure forever.
- (Loss of Control) Once made available, disclosed data will inevitably be used for purposes (not necessarily known at the time of the collection) beyond the purpose for which it was disclosed. Underlying this assumption is the premise that the mere existence of something is sufficient to tempt people to use it in whatever way they see fit to suit their needs and desires. The public and private sector will inevitably find new uses to improve the efficiency, security, or reach of their operations; foregoing opportunities can easily result in a loss of competitive edge. Law enforcement agencies will inevitably seek access to the data in the belief that it will help their investigative practices.
- (Linkability) Data disclosed in one transaction will inevitably be linked to data disclosed in other transactions (if not for reasons related to security then for

marketing, inventory management, or efficiency purposes), unless the cost of linking outweighs the benefits. With the trend or at least the capability of organizations to merge their databases at ever decreasing cost, it is naive to believe that linkable data that is submitted to different locations will remain unlinked.

To empower individuals to control their own data, PKIs must meet a number of basic privacy goals:

- Without the ability to remain anonymous, individuals have no control over their own privacy. Anonymity serves as the base case for privacy. In many situations, anonymity does more than serve privacy. If there is one thing that can be learned from the dramatic rise in identity fraud, it is that the use of person identifiers more often enables than prevents fraud. Disciplines such as treatments for medical conditions have long acknowledged that misuse can only be prevented through patient anonymity. The explosive rise of identity fraud in recent years illustrates that the same should hold true for most transaction mechanisms.
- Forcing individuals to use fully anonymous communication and transaction methods is almost as much an invasion on privacy as the other way around. Different individuals have disparate privacy preferences; surveys by Equifax and Louis Harris & Associates indicate that about 55% of people are privacy “pragmatists,” who are willing to trade personal data depending on a number of factors, including the benefits they will receive in return. Another important reason not to hardwire absolute anonymity into communication and transaction systems is that in many situations anonymity does not benefit anyone. In recognition of these facts, privacy-enhancing digital certification mechanisms should not make the property of anonymity invariant, but should enable each individual to decide for him or herself how much data to disclose in each transaction; this is called the *selective disclosure* paradigm.
- Anonymity of each transaction by itself is a necessary but insufficient condition to prevent linking of different transactions; any correlation that exists in the transcripts of any two protocol executions, for instance in the form of a public key, may be used to link them. Without unlinkability individuals cannot control how much data they actually disclose, since the aggregate information learned by linking different transactions will typically reveal much more than the data items that were willingly disclosed on each separate occasion. Without control over the degree of linkability, the paradigm of selective disclosure loses its power with each new disclosure. In particular, if a person is identified in a single transaction, then all his or her past and future transactions become traceable. *Unlinkability* is essential to prevent gradual erosion of privacy.

- In case a certificate holder authenticates his or her certificate showings by means of a digital signature, he or she leaves a permanent self-authenticating record that can be verified by anyone. This gives coercive powers to the receiver and anyone else who sees the signed statement. If transactions are untraceable, little harm may come from this, but there is always the possibility that signed data disclosed in one anonymous transaction can be linked to later transactions in which an identifier is revealed. For this reason, it is desirable that individuals can authenticate messages and attribute data in a manner that does not leave a self-authenticating record.
- Given the enormous security benefits offered by smartcard-based implementations, there is an urgent need to preserve all these privacy properties in this setting. In particular, smartcards should be unable to leak personal data and other attributes stored inside, and should be unable to learn any information from the outside world other than what their holders consent to. These properties should hold in the strongest possible sense, namely in the presence of CAs that have access to cryptographic backdoors and conspire with certificate verifiers.

In Chapters 2 to 6 we will develop cryptographic techniques that meet these privacy objectives and overcome the security and efficiency problems in Section 1.1. In the next section we give an overview of these techniques.

1.3 Outlook

1.3.1 Basic building blocks

Throughout the rest of the book the term “digital certificate” will always refer to the CA’s signature only; it does not include the public key or any information associated with it by the CA’s signature. This convention is not mainstream²⁶ but makes it easier to distinguish between various cryptographic objects.

We start in Chapter 2 with an overview of the cryptographic preliminaries needed to understand the material in the other chapters. Several new primitives will be introduced that play a fundamental role in the other chapters. In particular, we introduce two functions that are one-way and collision-intractable, and for both we design practical techniques for proving knowledge of an inverse and for constructing digital signatures. We also introduce a new kind of digital certificates.

In Chapters 3 and 4 we develop two basic building blocks:

- A certificate *issuing protocol* with the following properties:

²⁶Many publications consider a certificate to be the data structure comprised of the CA’s signature, the public key it certifies, and any information assigned to that public key.

- The receiver receives an unforgeable triple of the form (secret key, public key, certificate). The (secret key, public key) of the triple is a key pair for use by the receiver, while the certificate is digital signature of the issuer, made using its own secret key.
 - The receiver can ensure that the (public key, certificate) pair of the triple is fully blinded. (Consequently, at least part of the secret key is blinded as well, since the public key corresponds uniquely to the secret key.)
 - The receiver cannot blind a non-trivial blinding-invariant part of the secret key of the triple. In this blinding-invariant part, the issuer can encode an arbitrary number of attributes.
- A certificate *showing protocol* with the following properties:
 - To show a retrieved triple, the receiver discloses the (public key, certificate) pair and uses the secret key of the triple to authenticate a message. (This authentication serves at the very least to prevent replay.) The authentication mechanism allows the receiver to avoid leaving behind a self-authenticating record.
 - The authentication mechanism is such that the receiver not only authenticates the message, but also demonstrates a property of the attributes encoded into its certified key pair. The receiver has full control over which property is demonstrated: it can be any satisfiable proposition from proposition logic, where the atomic propositions are relations that are linear in the encoded attributes. Any other information about the attributes remains unconditionally hidden.

An automated negotiation mechanism such as that of OPS/P3P could be used to implement the negotiation process in the showing protocol.

The certificate showing protocol techniques will be developed in Chapter 3 and the issuing protocol techniques in Chapter 4. In Section 5.1 we will show how to seamlessly combine the issuing and showing protocol techniques, without adding complexity and without compromising security or privacy.

The new certificates function in much the same way as do cash, stamps, cinema tickets, subway tokens, and so on: anyone can establish the validity of certificates and the non-identity data they certify, but no more than just that. A “demographic” certificate, for instance, can certify its holder’s age, income, marital status, and residence, all neatly tied to one public key by means of a single digital signature of the certificate issuer. Because the attributes are encoded into the certificate applicant’s secret key, certificate holders can decide for themselves, depending on the circumstances, which attributes to disclose. This goes beyond the analogy of using a marking pen to cross out data fields on a paper-based certificate; for instance, the

holder of a demographic certificate can prove that he or she is either over 65 or under 18, without revealing which is the case or anything else. Furthermore, actions involving different certificates cannot be linked on any other basis than by what is explicitly disclosed.

The basic building blocks are highly practical. They can be based on the RSA assumption as well as on the Discrete Logarithm assumption, and admit elliptic curve implementations with short public keys. The communication and computation complexity of the issuing protocol are virtually independent of the number of attributes encoded into a certified key pair, and the showing protocol is almost as efficient as protocols that cannot provide selective disclosure.

1.3.2 Additional privacy techniques

Section 5.2 is devoted to additional techniques to improve privacy for certificate holders:

- (Anonymous updating) In many cases one's right to access a service comes from a pre-existing relationship in which identity has already been established. We provide a technique that enables an individual to anonymously present a certified public key for updating to the CA. The CA can recertify the attributes, or updated versions of them, without needing to know their current values. A special application is to prevent the CA from learning the entire set of attributes of a certificate applicant. Different CAs can even certify different attributes for the same certified key pair.
- (Simulatable certificate information) To prevent online certificate repositories from serving as data warehouses containing indisputable information about certificate holders, so-called secret-key certificates (developed in Section 2.6) may be used. These certificates allow anyone to generate directory entries that are indistinguishable from the entries that list certificates issued by the CA, yet offer the same basic security. Secret-key certificates also have the advantage that a showing protocol execution is entirely zero-knowledge when the attribute property is demonstrated in zero-knowledge.
- (Hiding participation in a PKI) Using secret-key certificates, users can simulate certified public keys for PKIs in which they do not or may not participate. They can prove to be a participant of (at least) one out of many PKIs or to have attributes certified by a subset of several CAs, without revealing more. This reduces the scope for discrimination on the basis of one's (lack of) PKI access rights.
- (Selective disclosure for multiple attribute certificates) Rather than encoding many attributes into a single certified key pair, it may be preferable to distribute them across multiple certified key pairs. This helps avoid the aggregation of

an individual's attributes by a single CA, improves efficiency, and removes the need to update certificates more frequently than otherwise needed. Our selective disclosure techniques can be applied not only to attributes encoded into a single certified key pair, but also to attributes in different key pairs (possibly certified by different CAs). Likewise, different certificate holders can jointly demonstrate that their combined attributes meet certain properties.

- (Self-linkability) Certificate holders can anonymously prove in a simple variation of the showing protocol to be the originator of a plurality of showing protocol executions. As a special application, we show how to enable certificate holders in the showing protocol to build up reputations with organizations.

In Section 5.3 we will describe techniques to improve the privacy of certificate verifiers. Specifically, we will show how to perform the showing protocol in such a manner that the verifier receives a signed statement that proves that a certificate has been shown but unconditionally hides all or part of the attribute property that has been demonstrated. In applications where verifiers submit their showing protocol transcripts to the CA, for instance to enable the CA to detect and combat fraud, this property prevents the CA from learning which formulae the verifiers require their customers to demonstrate. At the same time, verifiers are unable to provide false information to the CA.

Our use of certified public keys has two side benefits: a secure session can be established without enabling wiretappers to identify the session initiator from its certified public key, and fraudulent CAs cannot falsely revoke certified public keys that are used only once.

1.3.3 Security techniques

In Section 5.4 we will show how to combine our issuing and showing protocols in such a manner that either one of the following two properties is achieved:

- (Unlimited-show certificates) Even if a certificate is shown an arbitrary number of times, the information that is revealed is no more than the aggregate information that is willingly disclosed in each of the individual showing protocol executions. (Multiple showings of the same certificate are all linkable, though; a certified public key in effect is a digital pseudonym.)
- (Limited-show certificates) If and only if a certificate is shown more than a predetermined number of times, the aggregate information that is revealed allows the computation of the entire secret key of the certificate holder (and in particular all the encoded attributes). The threshold can be arbitrarily set.

The limited-show property holds even if the certificate holder is free to choose the attribute property that it demonstrates in each showing protocol execution, and can be

combined with the verifier privacy technique described in Section 5.3. (That is, the CA will be able to trace perpetrators regardless of whether certificate verifiers hide a part of the formulae demonstrated.) Even conspiring certificate holders and verifiers cannot defeat the limited-show property.

The limited-show technique is highly practical: to compute one of the hidden attributes (for instance an identity attribute) in case of fraud, even in a military-strength implementation a “footprint” of a mere 60 bytes must be stored per showing protocol transcript, regardless of the complexity of the formula demonstrated and the number of encoded attributes.

In Section 5.5 we will show how to apply the limited-show techniques to discourage unauthorized lending and copying of certificates, and the deliberate discarding of certificates that contain attributes that the certificate holder does not want to show. These security techniques do not require tamper-resistant devices for certificate holders, nor do they require online certificate validation. When issuing gender or age certificates for gaining access to Internet discussion groups or Web sites, for instance, the issuer can encode into each certified key pair not only the designated receiver’s gender or age, but also some information that the receiver would like to keep secret (such as his or her credit card information, redeemable electronic coins, or an account access key). While the certificate holder can hide this secret when showing the certificate (by using our selective disclosure techniques), the certificate cannot be shown without actually knowing the encoded secret; lending therefore requires the certificate holder to give away the secret.

Furthermore, we show in Section 5.5 how to achieve non-repudiation for limited-show certificates, to prevent the CA from framing certificate holders by falsely claiming that limited-show certificates have been shown too many times. The evidence of fraud can be obtained in the form of a self-signed confession, and can be made unconditionally convincing. A particularly surprising feat is that the non-repudiation techniques can be made to work even when certificate applicants are anonymous to the certificate issuer.

We also describe in Section 5.5 measures to protect against leakage and misuse of the CA’s secret key, including measures to cope with attackers with infinite computing power (to prevent PKI meltdown). Another technique described in Section 5.5 concerns digital bearer certificates; these hide or do not contain any attributes that can be uniquely traced or linked to one person or to a select group.

Our techniques are not complementary to the currently prevailing ideas about digital certificates and PKIs, but encompass them as a special case. By way of example we will show in Section 5.5.1 how to encapsulate X.509v3 certificates. The new techniques are beneficial in any authentication-based communication or transaction environment in which there is no strict need to identify certificate holders at each and every occasion. The only acceptable role for X.509 and other identity certificates in such environments is to facilitate registration in case certificate applicants must be identified to the CA, similar to the way in which drivers’ licenses and passports

are traditionally used to acquire a permit or some other kind of authentication proof; even for this purpose, however, our certificates can be used.

In none of the techniques in this book do certificate verifiers need tamper-resistant devices.

1.3.4 Smartcard integration

All our techniques can be applied not only in the setting of software-only devices, but also in a setting where certificate holders in addition to a software-only device hold a smartcard. In Chapter 6 we first describe the many shortcomings of smartcard-only implementations, and list the advantages of combining smartcards with user-controlled software-only computers. We then show how to securely lift the software-only techniques of the preceding chapters to the smartcard setting in such a manner that the following privacy properties are guaranteed:

- The smartcard cannot learn the (public key, certificate) pair of its holder's certified key pairs, and cannot learn any encoded attributes its holder desires to keep secret.
- The smartcard cannot learn the property that is demonstrated in the showing protocol. In particular, regardless of the complexity of the formula demonstrated and the number of attributes encoded into a certified key pair, the smartcard performs exactly the same protocol. The smartcard cannot even decide whether multiple invocations of its assistance are for the purpose of showing the same certificate or different certificates, and can be prevented from learning any information on the number of certificates issued to its holder.
- All possible data leakages by and to the smartcard are prevented. This includes not only leakages that can be detected, but also subliminal channels. Consequently, the verifier learns nothing beyond the status of the formula(e) demonstrated; it cannot even distinguish whether the certificate holder is assisted by a smartcard or uses merely a software-only device.
- The smartcard can be prevented from developing any data that is statistically correlated to data known to the outside world (in particular, to the CA and certificate verifiers), so that even the contents of a returned smartcard that has been adversely programmed cannot reveal any information about the communications and transactions conducted by its holder (other than an upper bound on the number of showing protocol executions).

In this manner, the task of each smartcard is reduced to the absolute minimum, namely to protect the most basic security interests of the certificate issuer and its holder. This is desirable not only in light of privacy, but also for efficiency and security.

Our techniques accommodate situations in which the smartcard's task is deliberately broadened, for the purpose of controlling to which parties a certificate can be shown, which properties may be demonstrated, and so on. For instance, it may be desirable for individuals that their smartcard can assist in the showing protocol only when the designated verifier provides an identifier; as will be shown in Section 6.4.4, this suffices to protect against extortion attacks conducted over networks. Also, in some high-risk PKIs law enforcement may need the ability to trace the past actions of a designated certificate holder (but only with that person's awareness). In general, the smartcard can be prevented from learning anything beyond what it is expressly supposed to learn in order to perform a well-defined task known to its holder. This suffices to accommodate any legitimate needs to reduce the attainable privacy level.

The certificate issuing and showing protocols for software-only devices in Chapters 3 to 5 are a self-contained subset of the smartcard-enhanced protocols. An important advantage of this architecture is that all the security protections of the software-only system apply in the (presumably hypothetical) case that the tamper-resistance of a large number of smartcards is compromised overnight. It also enables PKI implementations in which some certificate holders hold software-only devices and others use smartcards. In particular, a PKI can be introduced as a software-only system and migrate gradually to a smartcard-enhanced system as the demand rises for greater efficiency, functionality, and security.

The computation and storage requirements for the smartcard do not depend on the number of encoded attributes or the complexity of the demonstrated formulae. Our smartcard techniques can even be implemented using low-cost 8-bit smartcards with limited memory and no cryptographic coprocessor, as will be shown in Section 6.5.1. This minimizes the cost for all parties, and allows manufacturers to devote the bulk of smartcard logic to improved tamper-resistance measures.

Different PKIs can make use of the same smartcard without being able to interchange personal data (unless the card holder consents). Certificate applications can be built on top of widely available smartcards that provide only basic identification or signature functionality.

When limited-show short-lived certificates are issued, the need to rely on (timely) revocation greatly reduces, and so our smartcard techniques also help overcome the cost, efficiency, and privacy problems of off-line and online certificate revocation and validation mechanisms. (Revocation of encryption keys can be avoided by randomly generating one-time encryption keys afresh at the start of each authenticated session.)

In Section 6.5 we show how certificate holders can securely return to the CA any retrieved certificates that have not yet been shown, how to discourage certificate holders from using their certificates to help remote parties gain access to services for which they do not hold the proper certificates themselves, how to design secure bearer certificates with optimal privacy, and how to prevent organizations and other verifiers from discriminating against certificate holders who do not disclose their built-in identifiers.

1.3.5 Security and privacy guarantees

All the techniques can be based on the RSA assumption as well as on the Discrete Logarithm assumption, and admit elliptic curve implementations with short public keys. Many of the security aspects can rigorously be proved equivalent to the security of either one of these assumptions in the so-called random oracle model.

All the privacy properties are guaranteed in the strongest possible sense: even if all the verifiers, smartcards, and CAs conspire in an active attack, are given infinite computing power, and jointly establish secret information in an preparatory phase, they cannot learn more than what can be inferred from the assertions that are voluntarily demonstrated in executions of the showing protocol. Consequently, individuals can prevent secondary use of their attribute data and can at all times ensure the correctness, timeliness, and relevance of their own data. At the same time, the risk of identity fraud is minimized.

While this information-theoretical privacy guarantee is very strong, it is important to realize that computational privacy would be unsatisfactory:

- The infeasibility assumption at the heart of breaking computational privacy is based on a specific distribution of the system parameters and key material generated by the CA. It may be hard or even impossible to verify that these are indeed generated in accordance with the proper probability distribution. A clever method of generating the system parameters or the key material may enable the CA to trace communications and transactions with modest computational effort.
- Another danger of computational privacy is that one or two decades from now it may be entirely practical for CAs to retroactively trace any or all of today's communications and transactions. The expected advances in algorithmics and progression in sheer computing power²⁷ will make it possible then to break implementations based on key sizes deemed sufficiently strong today, without needing a polynomial-time attacking algorithm. Indeed, virtually all the cryptographic systems in use today employ keys that for efficiency reasons are as small as possible; these key sizes do not guarantee invulnerability for more than a decade.

In either of these two cases, the resulting level of privacy-intrusion may be much more damaging than that of PKIs without any form of privacy to begin with, because certificate holders will be less inhibited in their actions.

The difference between computational and information-theoretical privacy can be viewed as follows. With computational privacy all the secrets of individuals end up

²⁷In 1965, Moore predicted that the number of components on integrated circuits would double every year for ten years. In 1975 he predicted a doubling every two years instead of every year. Thus far, Moore's prediction has been remarkably accurate. It is anticipated that by the year 2010 we will be down to atomic dimensions.

in the outside world, encrypted under a public key. Information-theoretical privacy guarantees that the secrets do not get out there in the first place.

In a practical implementation of a design that offers information-theoretical privacy, certificate holders should be given the freedom to select their own method of generating the random numbers needed to protect their own privacy. Those who desire the strongest privacy level should use random bits produced by a noise generator (post-processed by arithmetical methods to remove correlations), while others may be comfortable using pseudorandom bit generators or other methods that offer at most computational privacy. The fundamental difference with a system that has computational privacy hardwired into its design is that each certificate holder is free to choose or produce his or her own source of randomness, including the security parameters and seed values for pseudorandom generators.

For a general discussion of the difference between privacy-protecting methods and methods that merely create the illusion of privacy, see the Epilogue.

Our techniques do not protect against wiretapping and traffic analysis, but allow the modular adoption of session encryption, anonymous remailers, and other measures as an additional layer. Techniques to prevent wiretapping and traffic analysis are largely platform dependent and not necessarily based on cryptography. Note also that confidentiality can be trivially achieved once the authenticity problem is solved; the authenticity proof can include a public key to be used for encryption. The independence of encryption is good design practice in any case, and avoids regulatory issues such as export controls.

1.3.6 Applicability

Our techniques facilitate a cookbook approach towards designing electronic communication and transaction systems. Applications of special interest include, but are not limited to: electronic cash; digital pseudonyms for public forums and virtual communities (such as Internet news groups and chat rooms); access control (to Virtual Private Networks, subscription-based services, Web sites, databases, buildings, and so on); digital copyright protection (anonymous certificates permitting use of works); electronic voting; electronic patient files; electronic postage; automated data bartering (integration with standardization efforts such as P3P is easy); online auctions; financial securities trading; pay-per-view tickets; public transport ticketing; electronic food stamps; road-toll pricing; national ID cards (but with privacy); permission-based marketing; Web site personalization; multi-agent systems; collaborative filtering (i.e., making recommendations to one person based on the opinions of like-minded persons); medical prescriptions; gift certificates; loyalty schemes; and, electronic gambling. The design of specific applications is outside the scope of this book, but is relatively straightforward in many cases.