

On Securing Untrusted Clouds with Cryptography

Yao Chen, Radu Sion

Stony Brook Network Security and Applied Cryptography Lab
{yaochen,sion}@cs.stonybrook.edu

Abstract

In a recent interview, Whitfield Diffie argued that “the whole point of cloud computing is economy” and while it is possible in principle for “computation to be done on encrypted data, [...] current techniques would more than undo the economy gained by the outsourcing and show little sign of becoming practical”. Here we explore whether this is truly the case and quantify just how expensive it is to secure computing in untrusted, potentially curious clouds.

We start by looking at the economics of computing in general and clouds in particular. Specifically, we derive the end-to-end cost of a CPU cycle in various environments and show that its cost lies between 0.5 picocents in efficient clouds and nearly 27 picocents for small enterprises (1 picocent = $\$1 \times 10^{-14}$), values validated against current cloud pricing.

We then explore the cost of common cryptography primitives as well as the viability of their deployment for cloud security purposes. We conclude that Diffie was correct. Securing outsourced data and computation against untrusted clouds is indeed costlier than the associated savings, with outsourcing mechanisms up to several orders of magnitudes costlier than their non-outsourced locally run alternatives.

1 Introduction

Commoditized outsourced computing has finally arrived, mainly due to the emergence of fast and cheap networking and efficient large scale computing. Amazon, Google, Microsoft and Oracle are just a few of the providers starting to offer increasingly complex storage and computation outsourcing. CPU cycles have become consumer merchandise.

In [10] we explored the end-to-end cost of a CPU cycle in various environments and show that its cost lies between 0.45 picocents in efficient clouds and 27 picocents for small business deployment scenarios (1 picocent = $\$1 \times 10^{-14}$). In terms of pure CPU cycle costs, current clouds present seemingly cost-effective propositions for personal and small enterprise clients.

Nevertheless, cloud clients are concerned with the **privacy of their data and computation** – this is often the primary adoption obstacle, especially for medium and large corporations, who often fall under strict regulatory compliance requirements. To address this, existing secure outsourcing research addressed several issues including guaranteeing integrity, confidentiality and privacy of outsourced data to secure querying on outsourced

Copyright 2012 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

Bulletin of the IEEE Computer Society Technical Committee on Data Engineering

encrypted database. Such assurances will likely require strong cryptography as part of elaborate intra- and client-cloud protocols. Yet, strong crypto is expensive. Thus, it is important to ask: how much cryptography can we afford in the cloud while maintaining the cost benefits of outsourcing?

Some believe the answer is simply *none*. For example, in a recent interview [48] Whitfield Diffie argued that “**current techniques would more than undo the economy [of] outsourcing and show little sign of becoming practical.**”

Here we set out to find out whether this holds and if so, by what margins. One way to look at this is in terms of CPU cycles. For each desired un-secured client CPU cycle, *how many additional cloud cycles can we spend on cryptography*, before its outsourcing becomes too expensive? We end up gaining the insight that today’s secure data outsourcing primitives are often orders of magnitude more expensive than local execution, mainly due to the fact that we do not know how to process complex functions on encrypted data efficiently enough. And outsourcing simple operations – such as existing research in querying encrypted data, keyword searches, selections, projections, and simple aggregates – is simply not profitable. Thus, while traditional security mechanisms allow the elegant handling of inter-client and outside adversaries, today it is still too costly to secure against cloud insiders with cryptography.

2 Cost Models

Parameters	H	S	M	L
CPU utilization	5-8%	10-12%	15-20%	40-56%
server:admin ratio	N.A.	100-140	140-200	800-1000
Space (sqft/month)	N.A.	\$0.5	\$0.5	\$0.25
PUE	N.A.	2-2.5	1.6-2	1.2-1.5

Figure 1: Sample key parameters.

using them is to help differentiate a set of key parameters (Figure 1).

2.1 Levels

Home Users (H). We include this scenario as a baseline for a simple home setup containing several computers. This could correspond to individuals with spare time to maintain a small set of computers, or a small home-based enterprise without staffing costs.

Small Enterprises (S). We consider here any scenario involving an infrastructure of up to 1000 servers run in-house in a commercial enterprise. The cost structure will start to feature most of the usual suspects, including commercial energy and network pricing, cooling, space leases, staffing etc. Small enterprises can not afford custom hardware, efficient power-distribution, and cooling or dedicated buildings among others. More importantly, in addition to power distribution inefficiencies, due to their nature, small enterprises cannot be run at high utilization as they would be usually under the incidence of business cycles and its associated peak loads.

Mid-size Enterprises (M). We consider here setups of up to 10,000 servers, run by a corporation, often in its own dedicated data center(s). Mid-size enterprises might have some clout and access to better service deals for network service as well as more efficient cooling and power distribution. They are not fully global, yet could feature several centers across one or two time zones, allowing increased independence from local load cycles as well as the ability to handle daily peaks better by shifting loads across timezones. All the above results ultimately in increased utilization (20-25% est.) and overall efficiency.

Large Enterprises/Clouds (L). Clouds and large enterprises run over 10,000 servers, cross multiple time-zones, often literally at a global level, with large data centers distributed across all continents and often in tens to hundreds of countries. For example Google has built a 30-acre site in Dalles, Oregon, next to a hydroelectric

dam providing cheap power. The site is composed of 34,000 square feet buildings [26]. Especially in cloud setups, high speed networks allow global-wide distribution and integration of load from thousands of individual points of load. This in turn flattens the 24-hour overall load curve and allows for efficient peak handling and comparably high utilization factors (50-60% est. [22]). Cloud providers run the most efficient infrastructures, and often are at the forefront of innovation. In one notorious instance, Google for example asked Intel for chips tolerating more heat, to allow for a few degrees increase in data center operating temperatures - which in turn increases cooling efficiency by whole percentage points [36]. Moreover, clouds have access to bulk-pricing for network service from large ISPs, often one order of magnitude cheaper than mid-size enterprises.

2.2 Factors

We now consider the cost factors that come into play across all of the above levels. These can be divided into a set of inter-dependent vectors, including: hardware (servers, networking gear), building (floor space leasing), energy (running hardware and cooling), service (administration, staffing, software maintenance), and network service. Other breakdown layouts of these factors are possible.

Server Hardware. Hardware costs include servers, racks, power equipment, network equipment, cooling equipment etc. We will discuss network equipment later.

We note that these costs drop with time, likely even by the time this goes to print. For example, while many of the current documented mid-size deployments use single or multi-CPU System-X blade servers at around \$1-2000 each [25], large data centers deploy custom setups at about \$3000 for 4 CPUs, near-future developments could yield important changes.¹ We will be conservative and empirically assume home PC prices of around \$750/CPU, small and mid-size enterprise costs of around \$1000/CPU (for 2 CPU blades) and cloud-level costs of no more than \$500/CPU.

Energy. Energy in data centers does not only include power, computing and networking hardware but the entire support infrastructure, including cooling, physical security, and overall facilities. A simple rough way to infer power costs is by estimating the Power Usage Efficiency (PUE) of the data center. The PUE is a metric defined by the GreenGrid Consortium to evaluate the energy efficiency of a data center [20] ($PUE = \text{Total Power Usage} / \text{IT Equipment Power Usage}$).

We will assume 1.2-1.5 PUE for large enterprises, 1.6-2 PUE for mid-size enterprises and 2-2.5 for small enterprises [38]. Costs of electricity are relatively uniform and documented [2].

Service. Evaluating the staffing requirements for data centers is an extremely complex endeavor as it involves a number of components such as software development and management, hardware repair, maintenance of cooling, building, network and power services.

Analytical approaches are challenged by the sparsity of available relevant supporting data sets.

We deployed a set of commonly accepted rule of thumb values that have been empirically developed and validate well [23]: the server to administrator ratio varies from 2:1 up to experimental 2500:1 values due to different degrees of automation and data management. In deployment, small to mid-size data centers feature a ratio of 100-140:1 whereas cloud level centers can go up to 1000:1 [19, 22].

Network Hardware. To allow for analysis of network intensive protocols, we chose to separate network transport service costs from the other factors of impact in the bottom line for CPU cycle. Specifically, while the internal network infrastructure costs will be factored in the data center costs, network service will not. We will estimate separately the cost of transferring a bit reliably to/from the data center intermediated by outside ISPs' networks. Internal network infrastructure costs can be estimated by evaluating the number of required switches and routers. The design of scalable large economy network topology with high inter-node bandwidth for data centers is an ever ongoing research problem [39]. We base our results on some of the latest state of the

¹In one documented instance, e.g., Amazon is working with Rackable Systems to deliver an under \$700 AMD-based 6 CPU board dubbed CEMS (Cooperative Expendable Micro-Slice Servers) V3.

$$\begin{aligned}
\text{CycleCost} &= \frac{\text{Server} + \text{Energy} + \text{Service} + \text{Network} + \text{Floor}}{\text{Total Cycles}} \\
&= \frac{\lambda_s \cdot N_s / \tau_s + (w_p \cdot \mu + w_i \cdot (1 - \mu)) \cdot PUE \cdot \lambda_e + \frac{N_s}{\alpha} \cdot \lambda_p + \lambda_w \cdot N_w / \tau_w + \lambda_f \cdot \frac{(w_p \cdot \mu + w_i \cdot (1 - \mu)) \cdot PUE}{\beta}}{\mu \cdot \nu \cdot N_s} \quad (1)
\end{aligned}$$

art research, deploying fat tree interconnect structures. Fat trees have been shown to offer significantly lower overall hardware costs with good overall connectivity factors.

Floor Space. Floor space costs vary wildly, by location and use. While small to mid-size enterprises usually have data centers near their location (thus sometimes incurring office-level pricing), large companies such as Google and Microsoft tend to build data centers on owned land, in less populated place where the per sqft price can be brought down much lower, often amortized to zero over time.

We also note that floor surface is directly related to power consumption and cooling with designs supporting anywhere from 40 to 250 watt/sqft [17]. Thus, the overall power requirements (driven by CPUs) impact directly the required space.

2.3 The Costs

We start by evaluating the amortized dollar cost of a CPU cycle in equation (1). See notations in Figure 2 and various setups’ parameters in Figure 1.

Symbol	Definition
N_s, N_w	number of servers, switches
α	administrator: server ratio
β	watt per sq ft
λ_s, λ_w	server, switch price
λ_p, λ_f	personnel, floor cost/sec
λ_e	electricity price/(watt·sec)
μ	CPU utilization
ν	CPU frequency
τ_s, τ_w	servers, switches lifespan (5 y.)
w_p, w_i	server power at peak, idle

Figure 2: Notations for (1).

Provider	Picocents
Amazon EC2	0.93 - 2.36
Google AppEngine	up to 2.31
Microsoft Azure	up to 1.96

Figure 3: Current pricings.

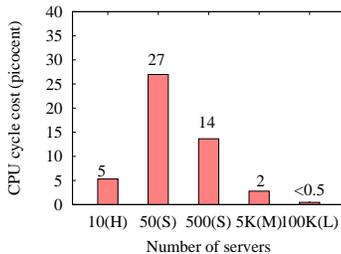


Figure 4: CPU cycle costs

“CPU cycles” are architecture-specific, yet we chose them to evade higher level, semantic-dependent units such as application-specific ‘transactions’. When reasoning about general computing it is not clear what types of higher-level transactions are appropriate to consider as ‘units’. After all, the dollar cost of serving an HTTP ‘transaction’ is only marginally relevant in evaluating the end-to-end costs of cloud-hosting arbitrary applications, across different infrastructures and languages. And, we will show that CPU cycles validate well as a consistent unit – probably in no small part due to the recent (almost) universality of x86 platforms across all environments, in effect reducing the impact of architecture specificity.

The results are depicted in Figure 4, costs ranging from 0.45 picocents/cycle in very large cloud settings all the way to (S), the costliest environment, where a cycle costs up to 27 picocents ($1 \text{ US picocent} = \1×10^{-14}). We validate our results by exploring the pricing of the main cloud providers (Figure 3). The prices lie surprisingly close to each other and to our estimates, ranging from 0.93 to 2.36 picocents/cycle. The difference in cost is due to the fact that these points include not only CPUs but also intra-cloud networking, instance-specific disk storage and cloud providers’ profit.

Storage Cost. Simply storing bits on disks has become truly cheap. Increased hardware reliability (with mean time between failures rated routinely above a million hours even for consumer markets) and economies of scale resulted in extreme drops in the costs of disks. In [10], we showed that in terms of amortized acquisition costs, the best price/hardware/MTBF ratio from our sample set is at 26.06 picocents/bit/year. The dominant factor is energy, 60-350 picocents/bit/year,

at 60-90% of the total cost. The lowest total cost from our sample set is at about 100 picocents/bit/year.

Network Service Published network service cost numbers place network service costs for large data centers

	AES-128	AES-192	AES-256
S	1.42E+03	1.48E+03	1.52E+03
L	2.37E+01	2.47E+01	2.53E+01

Figure 6: AES-128, AES-192, AES-256 costs (per byte) on 64-byte input.

	1024 bit		2048 bit	
	Encrypt	Decrypt	Encrypt	Decrypt
S	3.74E+06	1.03E+08	8.99E+06	6.44E+08
L	6.24E+04	1.72E+06	1.50E+05	1.07E+07

Figure 7: Cost of RSA encryption/decryption on 59-byte messages. (picocents)

at around \$13/ Mbps/ mo and for mid-size setups at \$95/Mbps/mo [22] for *guaranteed* bandwidth. Home user and small enterprise pricing benefits from economies of scale, e.g., Optimum Online provides 15/5 Mbps internet connection for small business starting at \$44.9/ mo [43]. Yet we note that the quoted bandwidth is not guaranteed and refers only to the hop connecting the client to the provider. Figure 5 summarizes network service cost in the four environments. When inferring the per-bit transmission costs we considered the uplink/downlink costs were independently priced at the same total price quoted for the entire connection. In other words, we assumed the provider would charge the same amount for only the uplink connection.

	H, S	M	L
monthly	\$44.90	\$95	\$13
bandwidth (d/u)	15/5 Mbps	per 1Mbps	per 1Mbps
dedicated	No	Yes	Yes
picocent/bit	115/345	3665	500

Figure 5: Summarized network service costs [10].

establishment, ACKs for sent data, window management, routing, packet parsing, re-transmissions). In the $S \rightarrow L$ scenario, it costs more than 900 picocents to transfer one bit reliably.

3 Cryptography

So far we know that a CPU cycle will set us back 0.45-27 picocents, transferring a bit costs at least 900 picocents, and storing it costs under 100 picocents/year. We now explore the costs of basic crypto and modular arithmetic. All values are in picocents. Note that CPU cycles needed in cryptographic operations often vary with optimization algorithms and types of hardware used (e.g., specialized secure CPUs and crypto accelerators with hardware RSA engines [1] are cheaper per cycle than general-purpose CPUs).

Symmetric Key Crypto. We first evaluate the per-bit costs of AES-128, AES-192, AES-256 and illustrate in Figure 6. The evaluation is based on results from the ECRYPT Benchmarking of Cryptographic Systems (eBACS) [6].

RSA. Using modular exponentiation, RSA public key encryption takes $O(k^2)$, private key decryption $O(k^3)$, and key generation $O(k^4)$ steps, where k is the number of bits in the modulus [28]. Numerous algorithms aim to improve the speed of RSA, mainly by reducing the time to do modular multiplications. In Figure 7, we illustrate the costs of RSA encryption/decryption using benchmark results from [6].

PK Signatures. We illustrate costs of DSA, and ECDSA signatures based on NIST elliptic curves [6] in Figures 8, 9.

Cryptographic Hashes We also show per byte cost of MD5 and SHA1 with varied input sizes.

4 Secure Outsourcing

Thus armed with an understanding of computation, storage, network and crypto costs, we now ask whether securing cloud computing against insiders is a viable endeavor.

	1024 bit		2048 bit	
	Sign	Verify	Sign	Verify
S	5.73E+07	6.94E+07	1.89E+08	2.30E+08
L	9.55E+05	1.16E+06	3.15E+06	3.84E+06

Figure 8: DSA on 59-byte messages. The 1024-bit DSA uses 148-byte secret key and 128-byte public key. The 2048-bit DSA uses 276-byte secret key and 256-byte public key.

	ECDSA-163		ECDSA-409	
	KG/SGN	Verify	KG/SGN	Verify
S	1.36E+08	2.65E+08	9.60E+08	1.91E+09
L	2.27E+06	4.41E+06	1.60E+07	3.19E+07
ECDSA-571				
	KG/SGN		Verify	
S	2.09E+09		4.18E+09	
L	3.48E+07		6.96E+07	

Figure 9: Costs of ECDSA signatures on 59-byte messages (curve over a field of size 2^{163} , 2^{409} , 2^{571} respectively). (picocents)

	MD5		SHA1	
	4096	64	4096	64
S	1.52E+02	3.75E+02	2.14E+02	6.44E+02
L	2.53E+00	6.25E+00	3.56E+00	1.07E+01

Figure 10: Per-byte cost of MD5 and SHA1 (with 64-byte and 4096-byte input).

We start by exploring what security means in this context. Naturally, the traditional usual suspects need to be handled in any outsourcing environment: (mutual) authentication, logic certification, inter-client isolation, network security as well as general physical security. Yet, all of these issues are addressed extensively in existing infrastructures and are not the subject of this work.

Similarly, for conciseness, within this scope, we will isolate the analysis from the additional costs of software patching, peak provisioning for reliability, network defenses etc.

4.1 Trust

We are concerned cloud clients being often reluctant to place sensitive data and logic onto remote servers without guarantees of compliance to their security policies [15, 29]. This is especially important in view of recent sub-poenas and other security incidents involving cloud-hosted data [11, 12, 35]. The viability of the cloud computing paradigm thus hinges directly on the issue of clients' trust and of major concern are cloud insiders. Yet how "trusted" are today's clouds from this perspective? We identify a set of scenarios.

Trusted clouds. In a *trusted* cloud, in the absence of unpredictable failures, clients are served correctly, in accordance to an agreed upon service contract and the cloud provider's policies. No insiders act maliciously.

Untrusted clouds. For *untrusted* clouds, we distinguish several cases depending on the types of illicit incentives existing for the cloud and the client policies with which these will directly conflict. We call a cloud *data-curious* if insiders thereof have incentives to violate confidentiality policies (mainly) for (sensitive) client data. Similarly, in an *access-curious* cloud, insiders will aim to infer client access patterns to data or reverse-engineer and understand outsourced computation logic. A *malicious* cloud will focus mainly on (data and computation) integrity policies and alter data or perform incorrect computation.

Reasonable cloud insiders are likely to factor in the potential illicit gains (the incentives to violate the policy), the penalty for getting caught, as well as the probability of detection. Thus for most practical scenarios,

insiders will engage in such behavior only if they can get away undetected with high probability, e.g., when no (cryptographic?) safeguards are in place to enable the detection.

4.2 Secure Outsourcing

Yet, millions of users embrace free web apps in **an untrusted provider model**. This shows that today's (mostly personal) cloud clients are willing to trade their privacy for (free) service. This is not necessarily a bad thing, especially at this critical-mass building stage, yet raises questions of clouds' viability for commercial, regulatory-compliant deployment, involving sensitive data and logic. And, from a bottom-line cost-perspective, is it worth even trying? This is what we aim to understand here.

In the following **we will assess whether clouds are economically tenable if their users do not trust them and therefore must employ cryptography and other mechanisms to protect their data**. A number of experimental systems and research efforts address the problem of outsourcing *data* to *untrusted service providers*, including issues ranging from searching in remote encrypted data to guaranteeing integrity and confidentiality to querying of outsourced data. In favor of cloud computing, we will set our analysis in the most favorable $S \rightarrow L$ scenario, which yields most CPU cycle savings.

4.2.1 The Case for Basic Outsourcing

Before we tackle cloud security, let us look at the simplest computation outsourcing scenario (where clients outsource data to the cloud, expect the cloud to process it, and send the results back). In existing work [10], we show that, to make (basic, unsecured) outsourcing cost effective, the cost savings (mainly from cheaper CPU cycles) need to outweigh the cloud's distance from clients. In $S \rightarrow L$, outsourced tasks should perform at least 1,000 CPU cycles per every 32 bit data, otherwise it is not worth outsourcing them.

4.2.2 Encrypted Data Storage with Integrity

With an understanding of the basic boundary condition defining the viability of outsourcing we now turn our attention to one of the most basic outsourcing scenarios in which a single data client places data remotely for simple storage purposes. In the $S \rightarrow L$ scenario, the amortized cost of storing a bit reliably either locally *or* *remotely* is under 9 picocents/month (including power). Network transfer however, is of at least 900 picocents per accessed bit, a cost that is not amortized and two orders of magnitude higher.

From a technological cost-centric point of view it is simply not effective to store data remotely: **outsourced storage costs can be upwards of 2+ orders of magnitude higher than local storage** for the $S \rightarrow L$ scenario *even in the absence of security assurances*.

Cost of Security. Yet, outsourced storage providers exist and thrive. This is likely due to factors outside of our scope, such as the convenience of being able to have access to the data from everywhere or collaborative application scenarios in which multiple data users share single data stores (multi-client settings). Notwithstanding the reason, since consumers have decided it is worth paying for outsourced storage, the next question we ask is, how much more would security cost in this context? We first survey some of the existing work.

Several existing systems encrypt data before storing it on potentially data-curious servers [7, 9, 37]. File systems such as I³FS [27], GFS [18], and Checksummed NCryptfs [46] perform online real-time integrity verification.

It can be seen that two main assurances are of concern here: integrity and confidentiality. The cheapest integrity constructs deployed in most of the above revolve around the use of hash-based MACs. As discussed above, SHA-1 based keyed MAC constructs with 4096-byte blocks would cost around 4 picocent/byte on the server and 200 picocents/byte on the client side, leading to a total cost of about 25 picocents/bit. This is at least 4 times lower than the cost of storing the bit for a year and at least one order of magnitude lower than

the costs incurred by transferring the same bit (at 900+ picocents/bit). Thus, **for outsourced storage, integrity assurance overheads are negligible.**

For publicly verifiable constructs, crypto-hash chains can help amortize their costs over multiple blocks. In the extreme case, a single signature could authenticate an entire file system, at the expense of increased I/O overheads for verification. Usually, a chain only includes a set of blocks.

For an average of twenty 4096 byte blocks² secured by a single hash-chain signed using 1024-bit RSA, would yield an amortized cost approximately 1M picocents per 4096-byte block (30+ picocents/bit) for client read verification and 180+ picocents/bit for write/signatures. This is up to **8 times more expensive than the MAC based case.**

4.2.3 Searches on Encrypted Data

Confidentiality alone can be achieved by encrypting the outsourced content before outsourcing to potentially access-curious servers. Once encrypted however, it cannot be easily processed by servers.

One of the first processing primitives that has been explored allows clients to search directly in remote encrypted data [4, 5, 13]. In these efforts, clients either linearly process the data using symmetric key encryption mechanisms, or, more often, outsource additional secure (meta)data mostly of size linear in the order of the original data set. This meta-data aids the server in searching through the encrypted data set while revealing as little as possible.

But is remote searching worth it vs. local storage? We concluded above that simply using a cloud as a remote file server is extremely non-profitable, up to several orders of magnitude. Could the searching application possibly make a difference? This would hold if either (i) the task of searching would be extremely CPU intensive allowing the cloud savings to kick in and offset the large losses due to network transfer, or (ii) the search is extremely selective and its results are a very small subset of the outsourced data set – thus amortizing the initial transfer cost over multiple searches.

We note that existing work does not support any complex search predicates outside of simple keyword matching search. Thus the only hope there is that the search-related CPU load (e.g., string comparison) will be enough cheaper in the cloud to offset the initial and result transfer costs.

Keyword searching can be done in asymptotically constant time, given enough storage or logarithmic if B-trees are used. While the client could maintain indexes and only deploy the cloud as a file server, we already discovered that this is not going to be profitable. Thus if we are to have any chance to benefit here, the index structures need to also be stored on the server.

In this case, the search cost includes the CPU cycle costs in reading the B-tree and performing binary searches within B-tree nodes. As an example, consider 32 bit search keys (e.g., as they can be read in one cycle from RAM), and a 1 TB database. 1-3 CPU cycles are needed to initiate the disk DMA per reading, and each comparison in the binary search requires another 1-3 cycles (for executing a comparison conditional jump operation). A B-tree with 16KB nodes will have approximately a 1000 fanout and a height of 4-5, so performing a search on this B-tree index requires about 100-300 CPU cycles. Thus in this simple remote search, $S \rightarrow L$ outsourcing would result in CPU-related savings of around 2,500-8,000 picocents per access. Transferring 32 bits from $S \rightarrow L$ costs upwards of 900 picocents. Outsourced searching becomes thus more expensive for any results upwards of 36 bytes per query.

4.2.4 Insights into Secure Query Processing

By now we start to suspect that similar insights hold also for outsourced query processing. This is because we now know that (i) the tasks to be outsourced should be CPU-intensive enough to offset the network overhead

²Douceur et al. [16], show that file sizes can be modeled using a log-normal distribution. E.g., for $\mu^e = 8.46$, $\sigma^e = 2.4$ and 20,000 files, the median file size would be 4KB, mean 80KB, along with a small number of files with sizes exceeding 1GB [3, 16].

– in other words, outsourcing peanut counting will never be profitable, and (ii) existing confidentiality (e.g., homomorphisms) and integrity (e.g., hash trees, aggregated signatures, hash chains) mechanisms can “secure” only very simple basic arithmetic (addition, multiplication) or data retrieval (selection, projection) which would cost under a few of cycles per word if done in an unsecured manner. In other words, *we do not know yet how to secure anything more complex than peanut counting*. And outsourcing of peanut counting is counter productive in the first place. Ergo our suspicion.

We start by surveying existing mechanisms. Hacigumus et al. [21] propose a method to execute SQL queries over partly obfuscated outsourced data to protect data **confidentiality** against a data-curious server. The main functionality relies on (i) partly obfuscating the outsourced data by dividing it into a set of partitions, (ii) query rewriting of original queries into querying referencing partitions instead of individual tuples, and (iii) client-side pruning of (necessarily coarse grained) results. The information leaked to the server is balancing a trade-off between client-side and server-side processing, as a function of the data segment size. [24] explores optimal bucket sizes for certain range queries.

Ge et al. [47] discuss executing aggregation queries with confidentiality on an untrusted server. Unfortunately, due to the use of extremely expensive homomorphisms this scheme leads to large processing times for any reasonably security parameter settings (e.g., for 1024 bit fields, 12+ days *per query* are required).

Other researchers have explored the issue of **correctness** in settings with potentially malicious servers. In a publisher-subscriber model, Devanbu et al. deployed Merkle trees to authenticate data published at a third party’s site [14], and then explored a general model for authenticating data structures [32, 33]. In [40, 41] as well as in [31], mechanisms for efficient integrity and origin authentication for selection predicate query results are introduced. Different signature schemes (DSA, RSA, Merkle trees [34] and BGLS [8]) are explored as potential alternatives for data authentication primitives. In [30, 45] *verification objects* VO are deployed to authenticate data retrieval in “edge computing” In [42, 44] Merkle tree and cryptographic hashing constructs are deployed to authenticate range query results.

To summarize, existing secure outsourced query mechanisms deploy (i) partitioning-based schemes and symmetric key encryption for (“statistical” only) confidentiality, (ii) homomorphisms for oblivious aggregation (SUM, COUNT) queries (simply too slow to be practical), (iii) hash trees/chains and (iv) signature chaining and aggregation to ensure correctness of selection/range queries and projection operators. SUM, COUNT, and projection usually behave linearly in the database size. Selection and range queries may be performed in constant time, logarithmic time or linear time depending on the queried attribute (e.g., whether it is a primary key) and the type of index used.

For illustration purposes, w.l.o.g., consider a scenario most favorable to outsourcing, i.e., assuming the operations behave linearly and are extremely selective, only incurring two 32-bit data transfers between the client and the cloud (one for the instruction and one for the result). Informally, to offset the network cost of $900 \times 32 \times 2 = 57,600$ picocents, only traversing a database of size at least 10^5 will generate enough CPU cycle cost savings. Thus it seems that with very selective queries (returning very little data) over large enough databases, outsourcing can break even.

Cost of Security. In the absence of security constructs, we were able to build a scenario for which outsourcing is viable. But what about a general scenario? What are the overheads of security there? It is important to understand whether the cost savings will be enough to offset them. While detailing individual secure query protocols is out of scope here, it is possible to reason generally and gain an insight into the associated order of magnitudes.

Existing integrity mechanisms deploy hash trees, hash chains and signatures to secure simple selection, projection or range queries. Security overheads would then include *at least* the (client-side) hash tree proof re-construction ($O(\log n)$ crypto-hashes) and subsequent signature verification of the tree’s root. The hash tree proofs are often used to authenticate range boundaries. The returned element set is then authenticated often through either a hash chain (in the case of range joins, at least 30 picocents per byte) or aggregated signature

constructs (e.g., roughly 60,000 picocents each, for selects or projections). This involves either modular arithmetic or crypto-hashing of the order of the result data set. For illustration purposes, we will again favor the case for outsourcing, and assume only crypto-hashing and a linear operation are applied.

Consider a database that has $n = 10^9$ tuples of 64 bits each. In that case (binary) hash tree nodes need to be at least 240 bits (80 + 160 bits = 2 pointers + hash value) long. If we assume 3 CPU cycles are needed per data item, the boundary condition results in selectivity $s \leq 0.00037$ before outsourcing starts to make economical sense. In a more typical scenario of $s = 0.001$ (queries are returning 0.1% of the tuples), a per-query loss of over 0.3 US cents will be incurred.

The above holds only for the $S \rightarrow L$ scenario in which hash trees are deployed. In the case of signature aggregation [41, 42], the break-even selectivity would be even lower due to the higher computation overheads.

5 To Conclude

In this paper we explored whether cryptography can be deployed to secure cloud computing against insiders. We estimated common cryptography costs (AES, MD5, SHA-1, RSA, DSA, and ECDSA) and finally explored outsourcing of data and computation to untrusted clouds. We showed that deploying the cloud as a simple remote encrypted file system is extremely unfeasible if considering only core technology costs. We also concluded that existing secure outsourced data query mechanisms are mostly cost-unfeasible because **today's cryptography simply lacks the expressive power to efficiently support outsourcing** to untrusted clouds. Hope is not lost however. We found borderline cases where outsourcing of simple range queries can break even when compared with local execution. These scenarios involve large amounts of outsourced data (e.g., 10^9 tuples) and extremely selective queries which return only an infinitesimal fraction of the original data (e.g., 0.00037%).

References

- [1] IBM 4764 PCI-X Cryptographic Coprocessor. Online at <http://www-03.ibm.com/security/cryptocards/pcixcc/overview.shtml>, 2007.
- [2] E. I. Administration. "average retail price of electricity to ultimate customers by end-use sector, by state". Online at http://www.eia.doe.gov/cneaf/electricity/epm/table5_6_a.html.
- [3] N. Agrawal, W. J. Bolosky, J. R. Douceur, and J. R. Lorch. A five-year study of file-system metadata. In *Proceedings of the 5th USENIX conference on File and Storage Technologies (FAST 07)*, Berkeley, CA, USA, 2007. USENIX Association.
- [4] G. Amanatidis, A. Boldyreva, and A. O'Neill. Provably-secure schemes for basic query support in outsourced databases. In S. Barker and G.-J. Ahn, editors, *DBSec*, volume 4602 of *Lecture Notes in Computer Science*, pages 14–30. Springer, 2007.
- [5] M. Bellare, A. Boldyreva, and A. O'Neill. Deterministic and efficiently searchable encryption. In A. Menezes, editor, *CRYPTO*, volume 4622 of *Lecture Notes in Computer Science*, pages 535–552. Springer, 2007.
- [6] D. J. Bernstein and T. L. (editors). ebacs: Ecrypt benchmarking of cryptographic systems. Online at <http://bench.cr.yp.to> accessed 30 Jan. 2009.
- [7] M. Blaze. A Cryptographic File System for Unix. In *Proceedings of the first ACM Conference on Computer and Communications Security*, pages 9–16, Fairfax, VA, 1993. ACM.
- [8] D. Boneh, C. Gentry, B. Lynn, and H. Shacham. Aggregate and verifiably encrypted signatures from bilinear maps. In *EuroCrypt*, 2003.
- [9] G. Cattaneo, L. Catuogno, A. D. Sorbo, and P. Persiano. The Design and Implementation of a Transparent Cryptographic Filesystem for UNIX. In *Proceedings of the Annual USENIX Technical Conference, FREENIX Track*, pages 245–252, Boston, MA, June 2001.

- [10] Y. Chen and R. Sion. To cloud or not to cloud?: musings on costs and viability. In *Proceedings of the 2nd ACM Symposium on Cloud Computing, SOCC '11*, pages 29:1–29:7, New York, NY, USA, 2011. ACM.
- [11] CNN. Feds seek Google records in porn probe. Online at <http://www.cnn.com>, Jan. 2006.
- [12] CNN. YouTube ordered to reveal its viewers. Online at <http://www.cnn.com>, July 2008.
- [13] R. Curtmola, J. Garay, S. Kamara, and R. Ostrovsky. Searchable symmetric encryption: improved definitions and efficient constructions. In *CCS '06: Proceedings of the 13th ACM conference on Computer and communications security*, pages 79–88, New York, NY, USA, 2006. ACM.
- [14] P. T. Devanbu, M. Gertz, C. Martel, and S. G. Stubblebine. Authentic third-party data publication. In *IFIP Workshop on Database Security*, pages 101–112, 2000.
- [15] Donna Bogatin. Google Apps data risks: Security vs. privacy. Online at <http://blogs.zdnet.com/micro-markets/?p=1021>, Feb. 2007.
- [16] J. R. Douceur and W. J. Bolosky. A large-scale study of file-system contents. In *Proceedings of the ACM SIGMETRICS international conference on Measurement and modeling of computer systems*, pages 59–70. ACM New York, NY, USA, 1999.
- [17] J. Fetzer. Internet data centers:end user & developer requirements. Online at <http://www.utilityeda.com/Summer2006/Mares.pdf>.
- [18] S. Ghemawat, H. Gobioff, and S. T. Leung. The Google File System. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP '03)*, pages 29–43, Bolton Landing, NY, October 2003. ACM SIGOPS.
- [19] A. Greenberg, J. Hamilton, D. A. Maltz, and P. Patel. The cost of a cloud: Research problems in data center networks. In *SIGCOM Computer Communications Review*, 2009.
- [20] T. G. Grid. Green grid metrics: Describing data center power efficiency. Online at http://www.thegreengrid.org/gg_content/Green_Grid_Metrics_WP.pdf.
- [21] H. Hacigumus, B. Iyer, C. Li, and S. Mehrotra. Executing SQL over encrypted data in the database-service-provider model. In *Proceedings of the ACM SIGMOD international conference on Management of data*, pages 216–227. ACM Press, 2002.
- [22] J. Hamilton. Internet-scale service efficiency. Large Scale Distributed Systems & Middleware (LADIS 2008), 2008.
- [23] J. Hamilton. On designing and deploying internet-scale services. Technical report, Windows Live Services Platform, Microsoft, 2008.
- [24] B. Hore, S. Mehrotra, and G. Tsudik. A privacy-preserving index for range queries. In *Proceedings of ACM SIGMOD*, 2004.
- [25] IBM. IBM blade servers. Online at <http://www-03.ibm.com/systems/bladecenter/hardware/servers/>.
- [26] S. H. John Markoff. Hiding in plain sight, google seeks more power. Online at <http://www.nytimes.com/2006/06/14/technology/14search.html>.
- [27] A. Kashyap, S. Patil, G. Sivathanu, and E. Zadok. I3FS: An In-Kernel Integrity Checker and Intrusion Detection File System. In *Proceedings of the 18th USENIX Large Installation System Administration Conference (LISA 2004)*, pages 69–79, Atlanta, GA, November 2004. USENIX Association.
- [28] R. Lab. How fast is the RSA algorithm? Online at <http://www.rsa.com/rsalabs/node.asp?id=2215>.
- [29] Larry Dignan. Will you trust Google with your data? Online at <http://blogs.zdnet.com/BTL/?p=4544>, Feb. 2007.
- [30] M. Atallah and C. YounSun and A. Kundu. Efficient Data Authentication in an Environment of Untrusted Third-Party Distributors. In *24th International Conference on Data Engineering ICDE*, pages 696–704, 2008.
- [31] Maithili Narasimha and Gene Tsudik. DSAC: integrity for outsourced databases with signature aggregation and chaining. Technical report, 2005.

- [32] C. Martel, G. Nuckolls, P. Devanbu, M. Gertz, A. Kwong, and S. Stubblebine. A general model for authenticated data structures. Technical report, 2001.
- [33] C. Martel, G. Nuckolls, P. Devanbu, M. Gertz, A. Kwong, and S. G. Stubblebine. A general model for authenticated data structures. *Algorithmica*, 39(1):21–41, 2004.
- [34] R. Merkle. Protocols for public key cryptosystems. In *IEEE Symposium on Research in Security and Privacy*, 1980.
- [35] J. Merritt. What google searches and data mining mean for you. Online at <http://www.talkleft.com/story/2006/01/25/692/74066>.
- [36] C. Metz. Google demanding intel’s hottest chips? Online at http://www.theregister.co.uk/2008/10/15/google_and_intel/.
- [37] Microsoft Research. Encrypting File System for Windows 2000. Technical report, Microsoft Corporation, July 1999. www.microsoft.com/windows2000/techinfo/howitworks/security/encrypt.asp.
- [38] R. Miller. Microsoft: PUE of 1.22 for data center containers. Online at <http://www.datacenterknowledge.com/archives/2008/10/20/microsoft-pue-of-122-for-data-center-containers/>.
- [39] A.-F. Mohammad, L. Alexander, and V. Amin. A scalable, commodity data center network architecture. *SIGCOMM Comput. Commun. Rev.*, 38(4):63–74, 2008.
- [40] E. Mykletun, M. Narasimha, and G. Tsudik. Authentication and integrity in outsourced databases. In *Proceedings of Network and Distributed System Security (NDSS)*, 2004.
- [41] E. Mykletun, M. Narasimha, and G. Tsudik. Signature bouquets: Immutability for aggregated/condensed signatures. In *Computer Security - ESORICS 2004*, volume 3193 of *Lecture Notes in Computer Science*, pages 160–176. Springer, 2004.
- [42] M. Narasimha and G. Tsudik. Authentication of Outsourced Databases using Signature Aggregation and Chaining. In *Proceedings of DASFAA*, 2006.
- [43] Optimum. Optimum online plans. Online at <http://www.buyoptimum.com>.
- [44] H. Pang, A. Jain, K. Ramamritham, and K.-L. Tan. Verifying Completeness of Relational Query Results in Data Publishing. In *Proceedings of ACM SIGMOD*, 2005.
- [45] H. Pang and K.-L. Tan. Authenticating query results in edge computing. In *ICDE '04: Proceedings of the 20th International Conference on Data Engineering*, page 560, Washington, DC, USA, 2004. IEEE Computer Society.
- [46] G. Sivathanu, C. P. Wright, and E. Zadok. Enhancing File System Integrity Through Checksums. Technical Report FSL-04-04, Computer Science Department, Stony Brook University, May 2004. www.fsl.cs.sunysb.edu/docs/nc-checksum-tr/nc-checksum.pdf.
- [47] Tingjian Ge and Stan Zdonik. Answering aggregation queries in a secure system model. In *VLDB '07: Proceedings of the 33rd international conference on Very large data bases*, pages 519–530. VLDB Endowment, 2007.
- [48] Whitfield Diffie. How Secure Is Cloud Computing? Online at <http://www.technologyreview.com/computing/23951/>, Nov. 2009.