

Picture credit: Rick Bowmer/AP



How to use
the new 65-megawatt
Bluffdale supercomputer:
a gentle introduction
to cryptanalysis

D. J. Bernstein
University of Illinois at Chicago &
Technische Universiteit Eindhoven

Disclaimers

Picture credit: Rick Bowmer/AP



How to use
the new 65-megawatt
Bluffdale supercomputer:
a gentle introduction
to cryptanalysis

D. J. Bernstein
University of Illinois at Chicago &
Technische Universiteit Eindhoven

Disclaimers

1. I don't work for NSA.

Picture credit: Rick Bowmer/AP



How to use
the new 65-megawatt
Bluffdale supercomputer:
a gentle introduction
to cryptanalysis

D. J. Bernstein
University of Illinois at Chicago &
Technische Universiteit Eindhoven

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.

Picture credit: Rick Bowmer/AP



How to use
the new 65-megawatt
Bluffdale supercomputer:
a gentle introduction
to cryptanalysis

D. J. Bernstein
University of Illinois at Chicago &
Technische Universiteit Eindhoven

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.

Picture credit: Rick Bowmer/AP



How to use
the new 65-megawatt
Bluffdale supercomputer:
a gentle introduction
to cryptanalysis

D. J. Bernstein
University of Illinois at Chicago &
Technische Universiteit Eindhoven

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.

Picture credit: Rick Bowmer/AP



How to use
the new 65-megawatt
Bluffdale supercomputer:
a gentle introduction
to cryptanalysis
D. J. Bernstein
University of Illinois at Chicago &
Technische Universiteit Eindhoven

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.
5. Also *assuming* use of traditional transistors+wires, probably with some optics; plus long-term storage. Quantum computing would require different analysis.

credit: Rick Bowmer/AP



use

65-megawatt

supercomputer:

introduction

analysis

ernstein

ty of Illinois at Chicago &

the Universiteit Eindhoven

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.
5. Also *assuming* use of traditional transistors+wires, probably with some optics; plus long-term storage. Quantum computing would require different analysis.

Cryptogr

My miss

protect e

against e

owmer/AP



vatt

puter:

ion

is at Chicago &

siteit Eindhoven

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.
5. Also *assuming* use of traditional transistors+wires, probably with some optics; plus long-term storage. Quantum computing would require different analysis.

Cryptographic cha

My mission: Crypt
protect every Inter
against espionage-



Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.
5. Also *assuming* use of traditional transistors+wires, probably with some optics; plus long-term storage. Quantum computing would require different analysis.

ago &
hoven

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.
5. Also *assuming* use of traditional transistors+wires, probably with some optics; plus long-term storage. Quantum computing would require different analysis.

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage.

Disclaimers

1. I don't work for NSA.
2. NSA hasn't told me anything.
3. This is not a leak.
4. I'm *assuming* that NSA is not stupid.
5. Also *assuming* use of traditional transistors+wires, probably with some optics; plus long-term storage.
Quantum computing would require different analysis.

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage.

User needs crypto to be fast on devices designed primarily for doing something else:



ers

't work for NSA.

hasn't told me anything.

is not a leak.

assuming that
not stupid.

assuming use of
al transistors+wires,
y with some optics;
g-term storage.
m computing would
different analysis.

Cryptographic challenges

My mission: Cryptographically
protect every Internet packet
against espionage+sabotage.

User needs crypto to be fast
on devices designed primarily
for doing something else:



User also
crypto to

Some ex

- 2009 e
signature
(small
- 2010 e
signature
(trivial
- 2012 e
signature
(some

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage.

User needs crypto to be fast on devices designed primarily for doing something else:



User also needs crypto to be **secure**

Some examples of

- 2009 exploit of signatures in TLS (small public com)
- 2010 exploit of signatures in PlayStation (trivial—stupid S)
- 2012 exploit of signatures by Flare (somewhat large)

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage.

User needs crypto to be fast on devices designed primarily for doing something else:



User also needs crypto to be **secure**.

Some examples of crypto fail

- 2009 exploit of RSA-512 signatures in TI calculator (small public computation)
- 2010 exploit of ECDSA signatures in PlayStation 3 (trivial—stupid Sony mistake)
- 2012 exploit of MD5-based signatures by Flame malware (somewhat larger computation)

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage.

User needs crypto to be fast on devices designed primarily for doing something else:



User also needs crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512 signatures in TI calculators (small public computation);
- 2010 exploit of ECDSA signatures in PlayStation 3 (trivial—stupid Sony mistake);
- 2012 exploit of MD5-based signatures by Flame malware (somewhat larger computation).

Cryptographic challenges

My mission: Cryptographically protect every Internet packet against espionage+sabotage.

User needs crypto to be fast on devices designed primarily for doing something else:



User also needs crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512 signatures in TI calculators (small public computation);
- 2010 exploit of ECDSA signatures in PlayStation 3 (trivial—stupid Sony mistake);
- 2012 exploit of MD5-based signatures by Flame malware (somewhat larger computation).

Presumably many more examples not known to the public.

Graphic challenges

ion: Cryptographically
every Internet packet
espionage+sabotage.

eds crypto to be fast
es designed primarily
g something else:



User also needs
crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512
signatures in TI calculators
(small public computation);
- 2010 exploit of ECDSA
signatures in PlayStation 3
(trivial—stupid Sony mistake);
- 2012 exploit of MD5-based
signatures by Flame malware
(somewhat larger computation).

Presumably many more examples
not known to the public.

Critical c

Which c
fit the u
⇒ optim
cryptosy
for each

Challenges

cryptographically
internet packet
+ sabotage.

to be fast
and primarily
nothing else:



User also needs
crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512
signatures in TI calculators
(small public computation);
- 2010 exploit of ECDSA
signatures in PlayStation 3
(trivial—stupid Sony mistake);
- 2012 exploit of MD5-based
signatures by Flame malware
(somewhat larger computation).

Presumably many more examples
not known to the public.

Critical questions:

Which cryptographic
fit the user's cost
⇒ optimize choice
cryptosystem + a
for each user device

User also needs
crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512
signatures in TI calculators
(small public computation);
- 2010 exploit of ECDSA
signatures in PlayStation 3
(trivial—stupid Sony mistake);
- 2012 exploit of MD5-based
signatures by Flame malware
(somewhat larger computation).

Presumably many more examples
not known to the public.

Critical questions:

Which cryptographic system
fit the user's cost constraint
⇒ optimize choice of
cryptosystem + algorithm
for each user device.

User also needs
crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512
signatures in TI calculators
(small public computation);
- 2010 exploit of ECDSA
signatures in PlayStation 3
(trivial—stupid Sony mistake);
- 2012 exploit of MD5-based
signatures by Flame malware
(somewhat larger computation).

Presumably many more examples
not known to the public.

Critical questions:

Which cryptographic systems
fit the user's cost constraints?

⇒ optimize choice of

cryptosystem + algorithm

for each user device.

User also needs
crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512
signatures in TI calculators
(small public computation);
- 2010 exploit of ECDSA
signatures in PlayStation 3
(trivial—stupid Sony mistake);
- 2012 exploit of MD5-based
signatures by Flame malware
(somewhat larger computation).

Presumably many more examples
not known to the public.

Critical questions:

Which cryptographic systems
fit the user's cost constraints?

⇒ optimize choice of
cryptosystem + algorithm
for each user device.

Which cryptographic systems
can be broken by attackers?

⇒ optimize choice of
attack algorithm + device
for each cryptosystem.

User also needs
crypto to be **secure**.

Some examples of crypto failing:

- 2009 exploit of RSA-512
signatures in TI calculators
(small public computation);
- 2010 exploit of ECDSA
signatures in PlayStation 3
(trivial—stupid Sony mistake);
- 2012 exploit of MD5-based
signatures by Flame malware
(somewhat larger computation).

Presumably many more examples
not known to the public.

Critical questions:

Which cryptographic systems
fit the user's cost constraints?

⇒ optimize choice of
cryptosystem + algorithm
for each user device.

Which cryptographic systems
can be broken by attackers?

⇒ optimize choice of
attack algorithm + device
for each cryptosystem.

Heavy interactions between
high-level algorithms and
low-level computer architecture.

o needs
o be **secure**.

Examples of crypto failing:

exploit of RSA-512

ures in TI calculators

(public computation);

exploit of ECDSA

ures in PlayStation 3

(—stupid Sony mistake);

exploit of MD5-based

ures by Flame malware

(what larger computation).

Probably many more examples

known to the public.

Critical questions:

Which cryptographic systems
fit the user's cost constraints?

⇒ optimize choice of

cryptosystem + algorithm

for each user device.

Which cryptographic systems
can be broken by attackers?

⇒ optimize choice of

attack algorithm + device

for each cryptosystem.

Heavy interactions between

high-level algorithms and

low-level computer architecture.

Theory v

Predictio

physicist

sometim

Common

underlyin

calculati

re.

crypto failing:

RSA-512

calculators

mputation);

ECDSA

ayStation 3

Sony mistake);

MD5-based

ame malware

er computation).

more examples

public.

Critical questions:

Which cryptographic systems
fit the user's cost constraints?

⇒ optimize choice of

cryptosystem + algorithm

for each user device.

Which cryptographic systems
can be broken by attackers?

⇒ optimize choice of

attack algorithm + device

for each cryptosystem.

Heavy interactions between
high-level algorithms and
low-level computer architecture.

Theory vs. experim

Predictions made
physicists are often
sometimes wrong.

Common sources of
underlying models
calculations from t

Critical questions:

Which cryptographic systems fit the user's cost constraints?

⇒ optimize choice of

cryptosystem + algorithm

for each user device.

Which cryptographic systems can be broken by attackers?

⇒ optimize choice of

attack algorithm + device

for each cryptosystem.

Heavy interactions between high-level algorithms and low-level computer architecture.

Theory vs. experiment

Predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error: underlying models of physics, calculations from those models.

Critical questions:

Which cryptographic systems fit the user's cost constraints?

⇒ optimize choice of

cryptosystem + algorithm

for each user device.

Which cryptographic systems can be broken by attackers?

⇒ optimize choice of

attack algorithm + device

for each cryptosystem.

Heavy interactions between high-level algorithms and low-level computer architecture.

Theory vs. experiment

Predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error:
underlying models of physics;
calculations from those models.

Critical questions:

Which cryptographic systems fit the user's cost constraints?

⇒ optimize choice of

cryptosystem + algorithm

for each user device.

Which cryptographic systems can be broken by attackers?

⇒ optimize choice of

attack algorithm + device

for each cryptosystem.

Heavy interactions between high-level algorithms and low-level computer architecture.

Theory vs. experiment

Predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error: underlying models of physics; calculations from those models.

Experiments aren't perfect but catch many errors; resolve many disputes; provide raw data leading to new theories; build more confidence than theory alone can ever produce.

questions:

cryptographic systems
user's cost constraints?

optimize choice of

system + algorithm

user device.

cryptographic systems
broken by attackers?

optimize choice of

algorithm + device

cryptosystem.

interactions between

algorithms and

computer architecture.

Theory vs. experiment

Predictions made by theoretical
physicists are often disputed,
sometimes wrong.

Common sources of error:
underlying models of physics;
calculations from those models.

Experiments aren't perfect
but catch many errors;
resolve many disputes;
provide raw data
leading to new theories;
build more confidence than
theory alone can ever produce.

Is physics

Of course

Every field

theoretical

regarding

experiment

measure

we compare

phic systems
constraints?
e of

algorithm

ce.

phic systems
attackers?

e of

+ **device**

tem.

s between

ms and

r architecture.

Theory vs. experiment

Predictions made by theoretical
physicists are often disputed,
sometimes wrong.

Common sources of error:
underlying models of physics;
calculations from those models.

Experiments aren't perfect
but catch many errors;
resolve many disputes;
provide raw data
leading to new theories;
build more confidence than
theory alone can ever produce.

Is physics uniquely
Of course not.

Every field of science
theoreticians make
regarding observab
experimental science
measure those phe
we compare the re

Theory vs. experiment

Predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error:
underlying models of physics;
calculations from those models.

Experiments aren't perfect
but catch many errors;
resolve many disputes;
provide raw data
leading to new theories;
build more confidence than
theory alone can ever produce.

Is physics uniquely error-proof?
Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

Theory vs. experiment

Predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error:
underlying models of physics;
calculations from those models.

Experiments aren't perfect
but catch many errors;
resolve many disputes;
provide raw data
leading to new theories;
build more confidence than
theory alone can ever produce.

Is physics uniquely error-prone?

Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

Theory vs. experiment

Predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error:
underlying models of physics;
calculations from those models.

Experiments aren't perfect
but catch many errors;
resolve many disputes;
provide raw data
leading to new theories;
build more confidence than
theory alone can ever produce.

Is physics uniquely error-prone?

Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

What if measurements are
too expensive to carry out?

Measurements start with
scaled-down experiments,
work up towards
the scale of interest.

vs. experiment

predictions made by theoretical physicists are often disputed, sometimes wrong.

Common sources of error:
- Building models of physics;
- Predictions from those models.

Measurements aren't perfect
- They have many errors;
- They lead to many disputes;
- They provide raw data
- They lead to new theories;
- They have more confidence than
- They are the only one that can ever produce.

Is physics uniquely error-prone?
Of course not.

Every field of science:
- Theoreticians make predictions regarding observable phenomena;
- Experimental scientists measure those phenomena;
- We compare the results.

What if measurements are too expensive to carry out?
Measurements start with scaled-down experiments, work up towards the scale of interest.

Algorithmic error-prone
Theoreticians regarding
These predictions are often disputed

ment

by theoretical

n disputed,

of error:

of physics;

those models.

t perfect

rors;

utes;

eories;

ence than

ever produce.

Is physics uniquely error-prone?

Of course not.

Every field of science:

theoreticians make predictions

regarding observable phenomena;

experimental scientists

measure those phenomena;

we compare the results.

What if measurements are

too expensive to carry out?

Measurements start with

scaled-down experiments,

work up towards

the scale of interest.

Algorithm analysis

error-prone field of

Theoreticians make

regarding algorithms

These predictions

disputed, sometimes

Is physics uniquely error-prone?
Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

What if measurements are
too expensive to carry out?
Measurements start with
scaled-down experiments,
work up towards
the scale of interest.

Algorithm analysis is another
error-prone field of science.

Theoreticians make predictions
regarding algorithm performance.
These predictions are often
disputed, sometimes wrong.

Is physics uniquely error-prone?
Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

What if measurements are
too expensive to carry out?
Measurements start with
scaled-down experiments,
work up towards
the scale of interest.

Algorithm analysis is another
error-prone field of science.

Theoreticians make predictions
regarding algorithm performance.
These predictions are often
disputed, sometimes wrong.

Is physics uniquely error-prone?
Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

What if measurements are
too expensive to carry out?
Measurements start with
scaled-down experiments,
work up towards
the scale of interest.

Algorithm analysis is another
error-prone field of science.

Theoreticians make predictions
regarding algorithm performance.
These predictions are often
disputed, sometimes wrong.

Particularly error-prone:
cryptanalytic extrapolations
from an academic computation
to a serious real-world attack.

Is physics uniquely error-prone?
Of course not.

Every field of science:
theoreticians make predictions
regarding observable phenomena;
experimental scientists
measure those phenomena;
we compare the results.

What if measurements are
too expensive to carry out?

Measurements start with
scaled-down experiments,
work up towards
the scale of interest.

Algorithm analysis is another
error-prone field of science.

Theoreticians make predictions
regarding algorithm performance.
These predictions are often
disputed, sometimes wrong.

Particularly error-prone:
cryptanalytic extrapolations
from an academic computation
to a serious real-world attack.

We catch errors, resolve disputes
by carrying out experiments:
actually running these algorithms
on the largest scale we can.

Is uniquely error-prone?
or not.

Field of science:

Theorists make predictions
regarding observable phenomena;
Experimental scientists
measure those phenomena;
They compare the results.

Measurements are
expensive to carry out?
Experiments start with
small experiments,
moving towards
the area of interest.

Algorithm analysis is another
error-prone field of science.

Theoreticians make predictions
regarding algorithm performance.
These predictions are often
disputed, sometimes wrong.

Particularly error-prone:
cryptanalytic extrapolations
from an academic computation
to a serious real-world attack.

We catch errors, resolve disputes
by carrying out experiments:
actually running these algorithms
on the largest scale we can.

1980s seen
“QS” factor
costs 2^{10}

error-prone?

nce:

e predictions

ble phenomena;

tists

enomena;

esults.

ments are

arry out?

rt with

iments,

st.

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone:
cryptanalytic extrapolations
from an academic computation
to a serious real-world attack.

We catch errors, resolve disputes
by carrying out experiments:
actually running these algorithms
on the largest scale we can.

1980s security eva

“QS” factorization

costs 2^{100} to break

ne?

ns

mena;

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone: cryptanalytic extrapolations from an academic computation to a serious real-world attack.

We catch errors, resolve disputes by carrying out experiments: actually running these algorithms on the largest scale we can.

1980s security evaluation:

“QS” factorization algorithm costs 2^{100} to break RSA-1024

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone: cryptanalytic extrapolations from an academic computation to a serious real-world attack.

We catch errors, resolve disputes by carrying out experiments: actually running these algorithms on the largest scale we can.

1980s security evaluation:

“QS” factorization algorithm costs 2^{100} to break RSA-1024.

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone: cryptanalytic extrapolations from an academic computation to a serious real-world attack.

We catch errors, resolve disputes by carrying out experiments: actually running these algorithms on the largest scale we can.

1980s security evaluation:

“QS” factorization algorithm costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone: cryptanalytic extrapolations from an academic computation to a serious real-world attack.

We catch errors, resolve disputes by carrying out experiments: actually running these algorithms on the largest scale we can.

1980s security evaluation:

“QS” factorization algorithm costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

1991 Adleman: NFS

won't beat QS for RSA-1024.

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone: cryptanalytic extrapolations from an academic computation to a serious real-world attack.

We catch errors, resolve disputes by carrying out experiments: actually running these algorithms on the largest scale we can.

1980s security evaluation:

“QS” factorization algorithm costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

1991 Adleman: NFS

won't beat QS for RSA-1024.

Subsequent experiments \Rightarrow

NFS is much faster; maybe 2^{80} ?

Algorithm analysis is another error-prone field of science.

Theoreticians make predictions regarding algorithm performance.

These predictions are often disputed, sometimes wrong.

Particularly error-prone: cryptanalytic extrapolations from an academic computation to a serious real-world attack.

We catch errors, resolve disputes by carrying out experiments: actually running these algorithms on the largest scale we can.

1980s security evaluation:

“QS” factorization algorithm costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

1991 Adleman: NFS won't beat QS for RSA-1024.

Subsequent experiments \Rightarrow NFS is much faster; maybe 2^{80} ?

Actual security of RSA-1024 is still a matter of dispute: e.g., 2009 Bos–Kaihara–Kleinjung–Lenstra–Montgomery oppose NIST's transition to RSA-2048.

m analysis is another
one field of science.
cians make predictions
g algorithm performance.
redictions are often
, sometimes wrong.
arly error-prone:
alytic extrapolations
academic computation
ous real-world attack.
h errors, resolve disputes
ing out experiments:
running these algorithms
argest scale we can.

1980s security evaluation:
“QS” factorization algorithm
costs 2^{100} to break RSA-1024.
1990 Pollard: new “NFS” .
1991 Adleman: NFS
won't beat QS for RSA-1024.
Subsequent experiments \Rightarrow
NFS is much faster; maybe 2^{80} ?
Actual security of RSA-1024 is
still a matter of dispute: e.g.,
2009 Bos–Kaihara–Kleinjung–
Lenstra–Montgomery oppose
NIST's transition to RSA-2048.

The attack
Enough
should re
on amount
required
But can
this amount

is another
of science.
e predictions
m performance.
are often
es wrong.
orone:
apolations
computation
world attack.
esolve disputes
periments:
hese algorithms
e we can.

1980s security evaluation:
“QS” factorization algorithm
costs 2^{100} to break RSA-1024.
1990 Pollard: new “NFS” .
1991 Adleman: NFS
won't beat QS for RSA-1024.
Subsequent experiments \Rightarrow
NFS is much faster; maybe 2^{80} ?
Actual security of RSA-1024 is
still a matter of dispute: e.g.,
2009 Bos–Kaihara–Kleinjung–
Lenstra–Montgomery oppose
NIST's transition to RSA-2048.

The attacker's sup
Enough theory+ex
should reach conse
on amount of com
required to break a
But can the attack
this amount of com

1980s security evaluation:

“QS” factorization algorithm
costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

1991 Adleman: NFS
won't beat QS for RSA-1024.

Subsequent experiments \Rightarrow
NFS is much faster; maybe 2^{80} ?

Actual security of RSA-1024 is
still a matter of dispute: e.g.,
2009 Bos–Kaihara–Kleinjung–
Lenstra–Montgomery oppose
NIST's transition to RSA-2048.

The attacker's supercomput

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation

1980s security evaluation:

“QS” factorization algorithm
costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

1991 Adleman: NFS

won't beat QS for RSA-1024.

Subsequent experiments \Rightarrow

NFS is much faster; maybe 2^{80} ?

Actual security of RSA-1024 is

still a matter of dispute: e.g.,

2009 Bos–Kaihara–Kleinjung–

Lenstra–Montgomery oppose

NIST's transition to RSA-2048.

The attacker's supercomputer

Enough theory+experiment
should reach consensus

on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

1980s security evaluation:

“QS” factorization algorithm
costs 2^{100} to break RSA-1024.

1990 Pollard: new “NFS”.

1991 Adleman: NFS

won't beat QS for RSA-1024.

Subsequent experiments \Rightarrow

NFS is much faster; maybe 2^{80} ?

Actual security of RSA-1024 is

still a matter of dispute: e.g.,

2009 Bos–Kaihara–Kleinjung–

Lenstra–Montgomery oppose

NIST's transition to RSA-2048.

The attacker's supercomputer

Enough theory+experiment
should reach consensus

on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Security evaluation:

Factorization algorithm
 2^{80} to break RSA-1024.

Illard: new "NFS".

Shor: NFS

Quantum for RSA-1024.

Recent experiments \Rightarrow
much faster; maybe 2^{80} ?

Security of RSA-1024 is

matter of dispute: e.g.,

Shor–Kaihara–Kleinjung–

Montgomery oppose

transition to RSA-2048.

The attacker's supercomputer

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Communit

Bill Dal

"Commu

more en

evaluation:
an algorithm
to break RSA-1024.
by “NFS”.
NFS
for RSA-1024.
arguments \Rightarrow
how long; maybe 2^{80} ?
for RSA-1024 is
dispute: e.g.,
–Kleinjung–
strongly oppose
to RSA-2048.

The attacker’s supercomputer

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Communication vs

Bill Dally, 2013.06
“Communication takes
more energy than

The attacker's supercomputer

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes
more energy than arithmetic

The attacker's supercomputer

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes
more energy than arithmetic”.

The attacker's supercomputer

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes
more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of
energy that we spend today
is on transferring data.”

The attacker's supercomputer

Enough theory+experiment
should reach consensus
on amount of computation
required to break a system.

But can the attacker perform
this amount of computation?

Hypothesize attacker resources.

This talk: \$2 billion, 65MW.

Alternative: millions of
compromised Internet computers.

The interesting part: analyze
optimal use of those resources.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes
more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of
energy that we spend today
is on transferring data.”

Depends what you're doing!

Computations fundamentally vary
in amount of communication
(distance and volume)
and amount of arithmetic.

Attacker's supercomputer

theory+experiment

reach consensus

amount of computation

to break a system.

the attacker perform

amount of computation?

resize attacker resources.

cost: \$2 billion, 65MW.

size: millions of

distributed Internet computers.

interesting part: analyze

use of those resources.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of

energy that we spend today

is on transferring data.”

Depends what you're doing!

Computations fundamentally vary

in amount of communication

(distance and volume)

and amount of arithmetic.

Some algorithms

Square root

n^2 arithmetic

Supercomputer

Experiment

Consensus

Computation

in a system.

Can we make

computations?

Can we use

more resources.

Examples of

supercomputers.

Example: analyze

large resources.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of energy that we spend today is on transferring data.”

Depends what you're doing!

Computations fundamentally vary in amount of communication (distance and volume) and amount of arithmetic.

Some algorithms use

Square matrix-vector

n^2 arithmetic.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic” .

Stephen S. Pawlowski,
2013.06.18: “The majority of energy that we spend today is on transferring data.”

Depends what you're doing!

Computations fundamentally vary in amount of communication (distance and volume) and amount of arithmetic.

Some algorithms using n^2 d

Square matrix-vector product
 n^2 arithmetic.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of energy that we spend today is on transferring data.”

Depends what you're doing!

Computations fundamentally vary in amount of communication (distance and volume) and amount of arithmetic.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of energy that we spend today is on transferring data.”

Depends what you're doing!

Computations fundamentally vary in amount of communication (distance and volume) and amount of arithmetic.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of energy that we spend today is on transferring data.”

Depends what you're doing!

Computations fundamentally vary in amount of communication (distance and volume) and amount of arithmetic.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Communication vs. arithmetic

Bill Dally, 2013.06.17:

“Communication takes more energy than arithmetic”.

Stephen S. Pawlowski,

2013.06.18: “The majority of energy that we spend today is on transferring data.”

Depends what you're doing!

Computations fundamentally vary in amount of communication (distance and volume) and amount of arithmetic.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Integrals in quantum chemistry,
many common iterations,
graph algorithms, etc.:
 n^4 arithmetic, sometimes more.

Communication vs. arithmetic

y, 2013.06.17:

unication takes

ergy than arithmetic” .

S. Pawlowski,

.18: “The majority of

hat we spend today

nsferring data.”

s what you’re doing!

ations fundamentally vary

nt of communication

e and volume)

ount of arithmetic.

Some algorithms using n^2 data:

Square matrix-vector product:

n^2 arithmetic.

FFT for input size n^2 :

$n^2 \lg n$ arithmetic.

Matrix-matrix product:

typically n^3 arithmetic

without Strassen etc.

Integrals in quantum chemistry,

many common iterations,

graph algorithms, etc.:

n^4 arithmetic, sometimes more.

Chip are

is enoug

all data

s. arithmetic

0.17:

takes

arithmetic” .

wski,

majority of

end today

data.”

’re doing!

damentally vary

munication

me)

ithmetic.

Some algorithms using n^2 data:

Square matrix-vector product:

n^2 arithmetic.

FFT for input size n^2 :

$n^2 \lg n$ arithmetic.

Matrix-matrix product:

typically n^3 arithmetic

without Strassen etc.

Integrals in quantum chemistry,

many common iterations,

graph algorithms, etc.:

n^4 arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$

is enough to store

all data for size- n^2

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Integrals in quantum chemistry,
many common iterations,
graph algorithms, etc.:
 n^4 arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$

is enough to store
all data for size- n^2 FFT.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Integrals in quantum chemistry,
many common iterations,
graph algorithms, etc.:
 n^4 arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$

is enough to store
all data for size- n^2 FFT.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Integrals in quantum chemistry,
many common iterations,
graph algorithms, etc.:
 n^4 arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$
is enough to store
all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$
is also enough for
 n^2 parallel ALUs.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Integrals in quantum chemistry,
many common iterations,
graph algorithms, etc.:
 n^4 arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$
is enough to store
all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$
is also enough for
 n^2 parallel ALUs.

FFT takes time n^ϵ ,
thanks to parallelism? No!
Routing the FFT data
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Some algorithms using n^2 data:

Square matrix-vector product:
 n^2 arithmetic.

FFT for input size n^2 :
 $n^2 \lg n$ arithmetic.

Matrix-matrix product:
typically n^3 arithmetic
without Strassen etc.

Integrals in quantum chemistry,
many common iterations,
graph algorithms, etc.:
 n^4 arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$

is enough to store
all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$

is also enough for
 n^2 parallel ALUs.

FFT takes time n^ϵ ,
thanks to parallelism? No!

Routing the FFT data
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

1981 Brent–Kung: need $n^{1+\epsilon}$
even without wire delays.

algorithms using n^2 data:

matrix-vector product:

arithmetic.

input size n^2 :

arithmetic.

matrix product:

n^3 arithmetic

Strassen etc.

s in quantum chemistry,

common iterations,

algorithms, etc.:

arithmetic, sometimes more.

Chip area $n^{2+\epsilon}$

is enough to store

all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$

is also enough for

n^2 parallel ALUs.

FFT takes time n^ϵ ,

thanks to parallelism? No!

Routing the FFT data

occupies area $n^{2+\epsilon}$

for time $n^{1+\epsilon}$.

1981 Brent-Kung: need $n^{1+\epsilon}$

even without wire delays.

Chip area

is enough

several n

Routing

occupies

for time

Typical

also occ

for time

Closer to

the ALU

although

using n^2 data:

tor product:

n^2 :

duct:

netic

etc.

um chemistry,

rations,

etc.:

ometimes more.

Chip area $n^{2+\epsilon}$

is enough to store
all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$

is also enough for
 n^2 parallel ALUs.

FFT takes time n^ϵ ,
thanks to parallelism? No!

Routing the FFT data
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

1981 Brent–Kung: need $n^{1+\epsilon}$
even without wire delays.

Chip area $n^{2+\epsilon}$

is enough to store
several $n \times n$ mat

Routing matrix pro
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 A
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dom
although not by m

ata:

Chip area $n^{2+\epsilon}$
is enough to store
all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$
is also enough for
 n^2 parallel ALUs.

FFT takes time n^ϵ ,
thanks to parallelism? No!
Routing the FFT data
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

1981 Brent–Kung: need $n^{1+\epsilon}$
even without wire delays.

ct:

Chip area $n^{2+\epsilon}$
is enough to store
several $n \times n$ matrices.

Routing matrix product
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 ALUs
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dominates,
although not by much.

stry,

ore.

Chip area $n^{2+\epsilon}$

is enough to store
all data for size- n^2 FFT.

Chip area $n^{2+\epsilon}$

is also enough for
 n^2 parallel ALUs.

FFT takes time n^ϵ ,

thanks to parallelism? No!

Routing the FFT data

occupies area $n^{2+\epsilon}$

for time $n^{1+\epsilon}$.

1981 Brent–Kung: need $n^{1+\epsilon}$

even without wire delays.

Chip area $n^{2+\epsilon}$

is enough to store
several $n \times n$ matrices.

Routing matrix product

occupies area $n^{2+\epsilon}$

for time $n^{1+\epsilon}$.

Typical n^3 arithmetic

also occupies n^2 ALUs

for time $n^{1+\epsilon}$.

Closer look at ϵ :

the ALU cost dominates,

although not by much.

area $n^{2+\epsilon}$
is enough to store
for size- n^2 FFT.

area $n^{2+\epsilon}$
is enough for
parallel ALUs.

requires time n^ϵ ,
no parallelism? No!
the FFT data
requires area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Ben-Kung: need $n^{1+\epsilon}$
without wire delays.

Chip area $n^{2+\epsilon}$
is enough to store
several $n \times n$ matrices.

Routing matrix product
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 ALUs
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dominates,
although not by much.

>90% of
of typical
is spent
<10% of
Is Bluffd

n^2 FFT.

ϵ ,
sm? No!

data

ϵ

need $n^{1+\epsilon}$
delays.

Chip area $n^{2+\epsilon}$
is enough to store
several $n \times n$ matrices.

Routing matrix product
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 ALUs
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dominates,
although not by much.

>90% of the cost
of typical supercom
is spent on commu
<10% on ALUs.

Is Bluffdale built t

Chip area $n^{2+\epsilon}$
is enough to store
several $n \times n$ matrices.

Routing matrix product
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 ALUs
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dominates,
although not by much.

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?

Chip area $n^{2+\epsilon}$
is enough to store
several $n \times n$ matrices.

Routing matrix product
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 ALUs
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dominates,
although not by much.

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?

Chip area $n^{2+\epsilon}$

is enough to store
several $n \times n$ matrices.

Routing matrix product
occupies area $n^{2+\epsilon}$
for time $n^{1+\epsilon}$.

Typical n^3 arithmetic
also occupies n^2 ALUs
for time $n^{1+\epsilon}$.

Closer look at ϵ :
the ALU cost dominates,
although not by much.

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?
No; NSA is not stupid.

Doubling number of ALUs
would cost <10% extra.
Would \approx double performance
of matrix-matrix product
and heavier-arith computations.

NSA's computations have a mix
of heavy arith and heavy comm.

a $n^{2+\epsilon}$
h to store
 $n \times n$ matrices.

matrix product
s area $n^{2+\epsilon}$
 $n^{1+\epsilon}$.

n^3 arithmetic
opies n^2 ALUs
 $n^{1+\epsilon}$.

ook at ϵ :
l cost dominates,
n not by much.

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?
No; NSA is not stupid.

Doubling number of ALUs
would cost <10% extra.
Would \approx double performance
of matrix-matrix product
and heavier-arith computations.

NSA's computations have a mix
of heavy arith and heavy comm.

GPUs ha
but relat
commun
a few lo
Is Bluffd

prices.

product

€

etic

ALUs

minates,

much.

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?

No; NSA is not stupid.

Doubling number of ALUs
would cost <10% extra.

Would \approx double performance
of matrix-matrix product
and heavier-arith computations.

NSA's computations have a mix
of heavy arith and heavy comm.

GPUs have many
but relatively little
communication ca
a few long wires to

Is Bluffdale built t

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?

No; NSA is not stupid.

Doubling number of ALUs
would cost <10% extra.

Would \approx double performance
of matrix-matrix product
and heavier-arith computations.

NSA's computations have a mix
of heavy arith and heavy comm.

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?

No; NSA is not stupid.

Doubling number of ALUs
would cost <10% extra.

Would \approx double performance
of matrix-matrix product
and heavier-arith computations.

NSA's computations have a mix
of heavy arith and heavy comm.

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?

>90% of the cost
of typical supercomputers
is spent on communication;
<10% on ALUs.

Is Bluffdale built this way?
No; NSA is not stupid.

Doubling number of ALUs
would cost <10% extra.

Would \approx double performance
of matrix-matrix product
and heavier-arith computations.

NSA's computations have a mix
of heavy arith and heavy comm.

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?
No; NSA is not stupid.

Adding communication
between adjacent ALUs
would cost very little.

Would drastically speed up
matrix-matrix product
and heavier-comm computations:
FFT, sorting, etc.

f the cost
al supercomputers
on communication;
n ALUs.

ale built this way?
A is not stupid.

g number of ALUs
ost <10% extra.
double performance
x-matrix product
vier-arith computations.

omputations have a mix
y arith and heavy comm.

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?
No; NSA is not stupid.

Adding communication
between adjacent ALUs
would cost very little.
Would drastically speed up
matrix-matrix product
and heavier-comm computations:
FFT, sorting, etc.

Docume
Intel Xeon
and a fe
plus adj
commun
Is Bluffd

computers
communication;

this way?
stupid.

of ALUs
extra.
performance
product
computations.

ns have a mix
heavy comm.

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?
No; NSA is not stupid.

Adding communication
between adjacent ALUs
would cost very little.
Would drastically speed up
matrix-matrix product
and heavier-comm computations:
FFT, sorting, etc.

Documentation te
Intel Xeon Phi has
and a few long wir
plus adjacent one-
communication (ri
Is Bluffdale built t

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?

No; NSA is not stupid.

Adding communication
between adjacent ALUs

would cost very little.

Would drastically speed up
matrix-matrix product
and heavier-comm computations:
FFT, sorting, etc.

Documentation tells me that
Intel Xeon Phi has many ALUs
and a few long wires to RAM
plus adjacent one-dimensional
communication (ring bus).

Is Bluffdale built this way?

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?

No; NSA is not stupid.

Adding communication

between adjacent ALUs

would cost very little.

Would drastically speed up

matrix-matrix product

and heavier-comm computations:

FFT, sorting, etc.

Documentation tells me that
Intel Xeon Phi has many ALUs
and a few long wires to RAM
plus adjacent one-dimensional
communication (ring bus).

Is Bluffdale built this way?

GPUs have many ALUs
but relatively little
communication capacity:
a few long wires to RAM.

Is Bluffdale built this way?

No; NSA is not stupid.

Adding communication

between adjacent ALUs

would cost very little.

Would drastically speed up

matrix-matrix product

and heavier-comm computations:

FFT, sorting, etc.

Documentation tells me that
Intel Xeon Phi has many ALUs
and a few long wires to RAM
plus adjacent one-dimensional
communication (ring bus).

Is Bluffdale built this way?

No; NSA is not stupid.

Adding **two-dimensional grid**

would drastically speed up

heavy-comm computations.

e.g. 1977 Thompson–Kung.

Grid examples: MasPar; FPGAs.

But FPGAs have other problems.

ave many ALUs
ctively little
nication capacity:
ng wires to RAM.
ale built this way?
A is not stupid.
communication
n adjacent ALUs
ost very little.
drastically speed up
matrix product
vier-comm computations:
orting, etc.

Documentation tells me that
Intel Xeon Phi has many ALUs
and a few long wires to RAM
plus adjacent one-dimensional
communication (ring bus).

Is Bluffdale built this way?
No; NSA is not stupid.

Adding **two-dimensional grid**
would drastically speed up
heavy-comm computations.
e.g. 1977 Thompson–Kung.

Grid examples: MasPar; FPGAs.
But FPGAs have other problems.

Save even
with 3D
e.g. 1983
Huge en
2D allow
energy in
up to ve
3D is ha
Some lin
(most in
presuma
Progress
e.g., $4 \times$
is often

ALUs

capacity:

to RAM.

this way?

stupid.

ation

t ALUs

tle.

speed up

duct

computations:

Documentation tells me that Intel Xeon Phi has many ALUs and a few long wires to RAM **plus** adjacent one-dimensional communication (ring bus).

Is Bluffdale built this way?

No; NSA is not stupid.

Adding **two-dimensional grid** would drastically speed up heavy-comm computations.
e.g. 1977 Thompson–Kung.

Grid examples: MasPar; FPGAs.

But FPGAs have other problems.

Save even more time with 3D arrangements

e.g. 1983 Rosenber

Huge engineering

2D allows easy scaling

energy input, heat

up to very large ch

3D is hard to scale

Some limited prog

(most interesting:

presumably used b

Progress often exa

e.g., $4 \times 16384 \times$

is often called “3D

Documentation tells me that Intel Xeon Phi has many ALUs and a few long wires to RAM **plus** adjacent one-dimensional communication (ring bus).

Is Bluffdale built this way?
No; NSA is not stupid.

Adding **two-dimensional grid** would drastically speed up heavy-comm computations.
e.g. 1977 Thompson–Kung.

Grid examples: MasPar; FPGAs.
But FPGAs have other problems.

Save even more time with 3D arrangement of ALU
e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of energy input, heat output up to very large chip area.
3D is hard to scale.

Some limited progress (most interesting: optics), presumably used by NSA.
Progress often exaggerated:
e.g., $4 \times 16384 \times 16384$ is often called “3D”.

tions:

Documentation tells me that Intel Xeon Phi has many ALUs and a few long wires to RAM **plus** adjacent one-dimensional communication (ring bus).

Is Bluffdale built this way?

No; NSA is not stupid.

Adding **two-dimensional grid**

would drastically speed up heavy-comm computations.

e.g. 1977 Thompson–Kung.

Grid examples: MasPar; FPGAs.

But FPGAs have other problems.

Save even more time with 3D arrangement of ALUs? e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of energy input, heat output up to very large chip area.

3D is hard to scale.

Some limited progress (most interesting: optics), presumably used by NSA.

Progress often exaggerated:

e.g., $4 \times 16384 \times 16384$

is often called “3D”.

ntation tells me that
on Phi has many ALUs
w long wires to RAM
acent one-dimensional
ication (ring bus).

ale built this way?
A is not stupid.

two-dimensional grid

astically speed up
omm computations.
7 Thompson–Kung.

mple: MasPar; FPGAs.
GAs have other problems.

Save even more time
with 3D arrangement of ALUs?
e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of
energy input, heat output
up to very large chip area.

3D is hard to scale.

Some limited progress
(most interesting: optics),
presumably used by NSA.

Progress often exaggerated:
e.g., $4 \times 16384 \times 16384$
is often called “3D”.

Special v

Typical c
between
better p
from AS

Some ex
ASICs b

calls me that
s many ALUs
res to RAM
-dimensional
ng bus).
his way?
upid.
nsional grid
speed up
utations.
son–Kung.
asPar; FPGAs.
other problems.

Save even more time
with 3D arrangement of ALUs?
e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of
energy input, heat output
up to very large chip area.

3D is hard to scale.

Some limited progress
(most interesting: optics),
presumably used by NSA.

Progress often exaggerated:
e.g., $4 \times 16384 \times 16384$
is often called “3D”.

Special vs. general

Typical cryptanaly
between $100\times$ and
better performance
from ASICs than f
mass-market CPU.

Some exceptions,
ASICs bring massi

Save even more time
with 3D arrangement of ALUs?
e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of
energy input, heat output
up to very large chip area.

3D is hard to scale.

Some limited progress
(most interesting: optics),
presumably used by NSA.

Progress often exaggerated:

e.g., $4 \times 16384 \times 16384$

is often called "3D".

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per trans
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup

Save even more time
with 3D arrangement of ALUs?
e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of
energy input, heat output
up to very large chip area.
3D is hard to scale.

Some limited progress
(most interesting: optics),
presumably used by NSA.
Progress often exaggerated:
e.g., $4 \times 16384 \times 16384$
is often called "3D".

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per transistor
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup.

Save even more time
with 3D arrangement of ALUs?
e.g. 1983 Rosenberg.

Huge engineering challenge.

2D allows easy scaling of
energy input, heat output
up to very large chip area.
3D is hard to scale.

Some limited progress
(most interesting: optics),
presumably used by NSA.
Progress often exaggerated:
e.g., $4 \times 16384 \times 16384$
is often called "3D".

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per transistor
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup.

Only in cryptanalysis? No.
Estimated ASIC improvement
from preliminary scan of other
supercomputing arith problems:
usually $>10\times$, often $>100\times$.

en more time

arrangement of ALUs?

3 Rosenberg.

engineering challenge.

vs easy scaling of

input, heat output

very large chip area.

rd to scale.

imited progress

interesting: optics),

bly used by NSA.

s often exaggerated:

: 16384×16384

called "3D".

Special vs. general purpose

Typical cryptanalytic arith:

between $100\times$ and $1000\times$

better performance per transistor

from ASICs than from

mass-market CPUs, GPUs.

Some exceptions, but overall

ASICs bring massive speedup.

Only in cryptanalysis? No.

Estimated ASIC improvement

from preliminary scan of other

supercomputing arith problems:

usually $>10\times$, often $>100\times$.

Frequent

chips sp

on decoo

\Rightarrow CPU,

reduce in

by adding

apply sa

to multi

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per transistor
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup.

Only in cryptanalysis? No.
Estimated ASIC improvement
from preliminary scan of other
supercomputing arith problems:
usually $>10\times$, often $>100\times$.

Frequent observations:
chips spend area, time
on decoding+scheduling
 \Rightarrow CPU/GPU designs
reduce insn-handling
by adding vectorization
apply same instructions
to multiple data/threads

Us?

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per transistor
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup.

Only in cryptanalysis? No.
Estimated ASIC improvement
from preliminary scan of other
supercomputing arith problems:
usually $>10\times$, often $>100\times$.

Frequent observation:
chips spend area, time, energy
on decoding+scheduling instr
 \Rightarrow CPU/GPU design trend:
reduce insn-handling cost
by adding vectorization—
apply same instruction
to multiple data/threads.

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per transistor
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup.

Only in cryptanalysis? No.
Estimated ASIC improvement
from preliminary scan of other
supercomputing arith problems:
usually $>10\times$, often $>100\times$.

Frequent observation:
chips spend area, time, energy
on decoding+scheduling insns.
 \Rightarrow CPU/GPU design trend:
reduce insn-handling cost
by adding vectorization—
apply same instruction
to multiple data/threads.

Special vs. general purpose

Typical cryptanalytic arith:
between $100\times$ and $1000\times$
better performance per transistor
from ASICs than from
mass-market CPUs, GPUs.

Some exceptions, but overall
ASICs bring massive speedup.

Only in cryptanalysis? No.
Estimated ASIC improvement
from preliminary scan of other
supercomputing arith problems:
usually $>10\times$, often $>100\times$.

Frequent observation:
chips spend area, time, energy
on decoding+scheduling insns.

\Rightarrow CPU/GPU design trend:
reduce insn-handling cost
by adding vectorization—
apply same instruction
to multiple data/threads.

But this does nothing
to reduce costs of
reading data from reg file,
writing data to reg file.

vs. general purpose

cryptanalytic arith:

100× and 1000×

performance per transistor

ICs than from

market CPUs, GPUs.

exceptions, but overall

bring massive speedup.

cryptanalysis? No.

ed ASIC improvement

eliminary scan of other

computing arith problems:

>10×, often >100×.

Frequent observation:

chips spend area, time, energy
on decoding+scheduling insns.

⇒ CPU/GPU design trend:

reduce insn-handling cost

by adding vectorization—

apply same instruction

to multiple data/threads.

But this does nothing

to reduce costs of

reading data from reg file,

writing data to reg file.

Obvious

reduce t

combine

doing m

between

Example

to comp

CPU rea

compute

With sep

CPU rea

writes; r

compute

purpose

tic arith:

d $1000\times$

e per transistor

from

s, GPUs.

but overall

ve speedup.

ysis? No.

mprovement

can of other

arith problems:

en $>100\times$.

Frequent observation:

chips spend area, time, energy
on decoding+scheduling insns.

\Rightarrow CPU/GPU design trend:

reduce insn-handling cost

by adding vectorization—

apply same instruction

to multiple data/threads.

But this does nothing

to reduce costs of

reading data from reg file,

writing data to reg file.

Obvious strategy to

reduce these reg c

combine arith ope

doing more arith

between read and

Example: Build ci

to compute $xy + z$

CPU reads regs $x,$

computes $xy + z$;

With separate mul

CPU reads x, y ; co

writes; reads back

computes $xy + z$;

Frequent observation:
chips spend area, time, energy
on decoding+scheduling insns.

⇒ CPU/GPU design trend:
reduce insn-handling cost
by adding vectorization—
apply same instruction
to multiple data/threads.

But this does nothing
to reduce costs of
reading data from reg file,
writing data to reg file.

Obvious strategy to
reduce these reg costs:
combine arith operations,
doing more arith
between read and write.

Example: Build circuit
to compute $xy + z$.
CPU reads regs x, y, z ;
computes $xy + z$; writes.

With separate mul, add:
CPU reads x, y ; computes x
writes; reads back; reads z ;
computes $xy + z$; writes.

Frequent observation:
chips spend area, time, energy
on decoding+scheduling insns.

⇒ CPU/GPU design trend:
reduce insn-handling cost
by adding vectorization—
apply same instruction
to multiple data/threads.

But this does nothing
to reduce costs of
reading data from reg file,
writing data to reg file.

Obvious strategy to
reduce these reg costs:
combine arith operations,
doing more arith
between read and write.

Example: Build circuit
to compute $xy + z$.

CPU reads regs x, y, z ;
computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ;
writes; reads back; reads z ;
computes $xy + z$; writes.

Observation:
end area, time, energy
ding+scheduling insns.

/GPU design trend:
insn-handling cost
g vectorization—
me instruction
ple data/threads.

does nothing
e costs of
data from reg file,
data to reg file.

Obvious strategy to
reduce these reg costs:
combine arith operations,
doing more arith
between read and write.

Example: Build circuit
to compute $xy + z$.

CPU reads regs x, y, z ;
computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ;
writes; reads back; reads z ;
computes $xy + z$; writes.

Common
evolved
Chip des
single-pr
eventual
circuit fo

ion:
time, energy
duling insns.

ign trend:

ng cost

ation—

ction

hreads.

ning

reg file,

g file.

Obvious strategy to
reduce these reg costs:
combine arith operations,
doing more arith
between read and write.

Example: Build circuit
to compute $xy + z$.

CPU reads regs x, y, z ;
computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ;
writes; reads back; reads z ;
computes $xy + z$; writes.

Common fp opera
evolved in this way

Chip designer saw
single-precision fp
eventually spent a
circuit for those m

gy
ns.

Obvious strategy to
reduce these reg costs:
combine arith operations,
doing more arith
between read and write.

Example: Build circuit
to compute $xy + z$.

CPU reads regs x, y, z ;
computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ;
writes; reads back; reads z ;
computes $xy + z$; writes.

Common fp operations
evolved in this way.

Chip designer saw many
single-precision fp muls,
eventually spent area on
circuit for those muls.

Obvious strategy to reduce these reg costs: combine arith operations, doing more arith between read and write.

Example: Build circuit to compute $xy + z$.

CPU reads regs x, y, z ; computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ; writes; reads back; reads z ; computes $xy + z$; writes.

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Obvious strategy to reduce these reg costs: combine arith operations, doing more arith between read and write.

Example: Build circuit to compute $xy + z$.

CPU reads regs x, y, z ; computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ; writes; reads back; reads z ; computes $xy + z$; writes.

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

Obvious strategy to reduce these reg costs: combine arith operations, doing more arith between read and write.

Example: Build circuit to compute $xy + z$.

CPU reads regs x, y, z ; computes $xy + z$; writes.

With separate mul, add:

CPU reads x, y ; computes xy ; writes; reads back; reads z ; computes $xy + z$; writes.

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

But people still run many single-precision computations. The multiplier transistors are **mostly sitting idle**.

strategy to
these reg costs:
arith operations,
ore arith
read and write.
e: Build circuit
ute $xy + z$.
nds regs x, y, z ;
es $xy + z$; writes.
parate mul, add:
nds x, y ; computes xy ;
eads back; reads z ;
es $xy + z$; writes.

Common fp operations
evolved in this way.

Chip designer saw many
single-precision fp muls,
eventually spent area on
circuit for those muls.

Then spent much more area
to expand the multiplier
to double-precision fp.

But people still run many
single-precision computations.
The multiplier transistors are
mostly sitting idle.

Another
Your app
mul-sub-
in its inn
Should C
include r
4 separa

to
osts:
rations,
write.
rcuit
z.
y, z;
writes.
l, add:
computes xy ;
; reads z ;
writes.

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

But people still run many single-precision computations. The multiplier transistors are **mostly sitting idle.**

Another example:
Your application does
mul-sub-sub-sub-s
in its inner loop.

Should CPU design
include mul circuit
4 separate sub circ

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

But people still run many single-precision computations. The multiplier transistors are **mostly sitting idle.**

Another example:
Your application does mul-sub-sub-sub-sub in its inner loop.

Should CPU designer include mul circuit, 4 separate sub circuits?

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

But people still run many single-precision computations. The multiplier transistors are **mostly sitting idle.**

Another example:

Your application does mul-sub-sub-sub-sub in its inner loop.

Should CPU designer include mul circuit, 4 separate sub circuits?

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

But people still run many single-precision computations.

The multiplier transistors are **mostly sitting idle.**

Another example:

Your application does mul-sub-sub-sub-sub in its inner loop.

Should CPU designer include mul circuit, 4 separate sub circuits?

Same CPU then runs another application.

Subtraction circuits are **mostly sitting idle.**

Common fp operations evolved in this way.

Chip designer saw many single-precision fp muls, eventually spent area on circuit for those muls.

Then spent much more area to expand the multiplier to double-precision fp.

But people still run many single-precision computations.

The multiplier transistors are **mostly sitting idle**.

Another example:

Your application does mul-sub-sub-sub-sub in its inner loop.

Should CPU designer include mul circuit, 4 separate sub circuits?

Same CPU then runs another application.

Subtraction circuits are **mostly sitting idle**.

CPU designer says no, reduces area per core.

⇒ Your application runs slowly.

n fp operations
in this way.

designer saw many
precision fp muls,
ly spent area on
or those muls.

ent much more area
nd the multiplier
e-precision fp.

ple still run many
recision computations.
Multiplier transistors are
sitting idle.

Another example:
Your application does
mul-sub-sub-sub-sub
in its inner loop.

Should CPU designer
include mul circuit,
4 separate sub circuits?

Same CPU then runs
another application.
Subtraction circuits are
mostly sitting idle.

CPU designer says no,
reduces area per core.

⇒ Your application runs slowly.

Many AS
beyond t
• Squari
than m
• Skip m
• Reduc
what i
• Add v
if appl
• etc.

tions
y.
many
muls,
rea on
muls.
more area
multiplier
n fp.
n many
mputations.
nsistors are
e.

Another example:
Your application does
mul-sub-sub-sub-sub
in its inner loop.

Should CPU designer
include mul circuit,
4 separate sub circuits?

Same CPU then runs
another application.
Subtraction circuits are
mostly sitting idle.

CPU designer says no,
reduces area per core.
⇒ Your application runs slowly.

Many ASIC fp spe
beyond today's CF

- Squaring is cheaper
than multiplication
- Skip most normal
- Reduce precision
what is actually
- Add very fast sq
if application ne
- etc.

Another example:
Your application does
mul-sub-sub-sub-sub
in its inner loop.

Should CPU designer
include mul circuit,
4 separate sub circuits?

Same CPU then runs
another application.
Subtraction circuits are
mostly sitting idle.

CPU designer says no,
reduces area per core.
⇒ Your application runs slowly.

Many ASIC fp speedups
beyond today's CPUs/GPUs

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Another example:

Your application does
mul-sub-sub-sub-sub
in its inner loop.

Should CPU designer
include mul circuit,
4 separate sub circuits?

Same CPU then runs
another application.

Subtraction circuits are
mostly sitting idle.

CPU designer says no,
reduces area per core.

⇒ Your application runs slowly.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Another example:

Your application does
mul-sub-sub-sub-sub
in its inner loop.

Should CPU designer
include mul circuit,
4 separate sub circuits?

Same CPU then runs
another application.

Subtraction circuits are
mostly sitting idle.

CPU designer says no,
reduces area per core.

⇒ Your application runs slowly.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Cryptanalysis involves
many multiplications
but also a much wider
variety of operations.

Even larger ASIC speedups.

example:
application does
-sub-sub-sub
inner loop.

CPU designer
mul circuit,
te sub circuits?

PU then runs
application.
tion circuits are
sitting idle.

signer says no,
area per core.
application runs slowly.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Cryptanalysis involves many multiplications but also a much wider variety of operations.
Even larger ASIC speedups.

So NSA
for each

The sma
ASIC de
Not a se
for \$2 bi

oes
ub
ner
t,
uits?
uns
n.
ts are
e.
s no,
ore.
n runs slowly.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Cryptanalysis involves many multiplications but also a much wider variety of operations.
Even larger ASIC speedups.

So NSA builds AS
for each application

The small problem
ASIC design effort
Not a serious issue
for \$2 billion.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Cryptanalysis involves many multiplications but also a much wider variety of operations.
Even larger ASIC speedups.

So NSA builds ASICs for each application?

The small problem:
ASIC design effort.
Not a serious issue for \$2 billion.

wly.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Cryptanalysis involves many multiplications but also a much wider variety of operations.

Even larger ASIC speedups.

So NSA builds ASICs for each application?

The small problem:
ASIC design effort.
Not a serious issue for \$2 billion.

Many ASIC fp speedups
beyond today's CPUs/GPUs:

- Squaring is cheaper than multiplication.
- Skip most normalizations.
- Reduce precision to what is actually needed.
- Add very fast sqrt if application needs it.
- etc.

Cryptanalysis involves many multiplications but also a much wider variety of operations.

Even larger ASIC speedups.

So NSA builds ASICs for each application?

The small problem:
ASIC design effort.
Not a serious issue for \$2 billion.

The big problem:
Unpredictable application mix.
NSA will want some agility to adapt to new computations and stop old computations.

Quantify using historical data:
how long is an ASIC useful?

ASIC fp speedups
today's CPUs/GPUs:
ing is cheaper
multiplication.
most normalizations.
e precision to
s actually needed.
ery fast sqrt
ication needs it.

alysis involves
multiplications
a much wider
of operations.
ger ASIC speedups.

So NSA builds ASICs
for each application?

The small problem:
ASIC design effort.
Not a serious issue
for \$2 billion.

The big problem:
Unpredictable application mix.
NSA will want some agility
to adapt to new computations
and stop old computations.

Quantify using historical data:
how long is an ASIC useful?

Obvious
some AS
mix of a
integrated
Take a g
Add exa
XYZZY
plus som
Think ab
XYZZZ?
Still sim
New CP
Merge s
if not m

speedups
CPUs/GPUs:
per
ion.
realizations.
n to
needed.
rt
eds it.
lves
ons
vider
ns.
speedups.

So NSA builds ASICs
for each application?

The small problem:
ASIC design effort.
Not a serious issue
for \$2 billion.

The big problem:
Unpredictable application mix.
NSA will want some agility
to adapt to new computations
and stop old computations.

Quantify using historical data:
how long is an ASIC useful?

Obvious solution for
some ASICs, plus
mix of **application**
integrated circuits

Take a general-purpose
Add exactly the bits
XYZZY needed by
plus some vectorization
Think ahead, add
XYZZZ? XZZZY? XYZZY?
Still similar cost to

New CPU for each
Merge similar applications
if not much cost in

So NSA builds ASICs
for each application?

The small problem:
ASIC design effort.

Not a serious issue
for \$2 billion.

The big problem:
Unpredictable application mix.
NSA will want some agility
to adapt to new computations
and stop old computations.

Quantify using historical data:
how long is an ASIC useful?

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU
Add exactly the big insns
XYZZY needed by application
plus some vectorization.

Think ahead, add agility:
XYZZZ? XZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application
Merge similar applications
if not much cost in area.

So NSA builds ASICs
for each application?

The small problem:
ASIC design effort.

Not a serious issue
for \$2 billion.

The big problem:

Unpredictable application mix.
NSA will want some agility
to adapt to new computations
and stop old computations.

Quantify using historical data:
how long is an ASIC useful?

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.

Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

builds ASICs
application?

all problem:
sign effort.

erious issue
illion.

problem:

stable application mix.

I want some agility
to new computations
old computations.

y using historical data:
g is an ASIC useful?

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.
Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide B

Critical t
and imp

ICs

on?

n:

.

e

lication mix.

ne agility

omputations

utations.

torical data:

IC useful?

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.

Add exactly the big insns
XYZZY needed by application,
plus some vectorization.

Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?

Still similar cost to ASIC.

New CPU for each application.

Merge similar applications
if not much cost in area.

1-slide Bluffdale us

Critical for algorithm
and implementor:

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

ix.
ns
ca:
Take a general-purpose CPU.
Add exactly the big insn
XYZZY needed by application,
plus some vectorization.
Think ahead, add agility:
XYZZZ? XZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm design
and implementor:

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.

Think ahead, add agility:
XYZZZ? XZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.

Think ahead, add agility:
XYZZZ? XZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Massive parallelism.

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.

Think ahead, add agility:
XYZZZ? XZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Massive parallelism.

Grid communication.

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.
Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Massive parallelism.

Grid communication.

Multiple instruction sets
with very useful instructions.

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.
Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Massive parallelism.

Grid communication.

Multiple instruction sets
with very useful instructions.

Some vectorization.

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.
Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Massive parallelism.

Grid communication.

Multiple instruction sets
with very useful instructions.

Some vectorization.

Occasional faults.

Obvious solution for NSA:
some ASICs, plus heterogeneous
mix of **application-tuned**
integrated circuits (ATICs).

Take a general-purpose CPU.
Add exactly the big insns
XYZZY needed by application,
plus some vectorization.
Think ahead, add agility:
XYZZZ? XZZZY? XYQZZY?
Still similar cost to ASIC.

New CPU for each application.
Merge similar applications
if not much cost in area.

1-slide Bluffdale user guide

Critical for algorithm designer
and implementor:

Massive parallelism.

Grid communication.

Multiple instruction sets
with very useful instructions.

Some vectorization.

Occasional faults.

Need to understand cryptanalysis:
ECM, sparse linear algebra,
differentials, FFTs, much more.