# Compiler issues: dependence analysis, vectorisation, automatic parallelisation

February 2018

Paul H J Kelly

The material for this part of the course is introduced only very briefly in Hennessy and Patterson. A good textbook which covers it properly is

- Michael Wolfe. High Performance Compilers for Parallel Computing. Addison Wesley, 1996.

Much of the presentation is taken from the following research paper:

- U. Banerjee. Unimodular transformations of double loops. In Proceedings of the Third Workshop on Programming Languages and Compilers for Parallel Computing, Irvine, CA. Pitman/MIT Press, 1990.

Banerjee's paper gives a simplified account of the theory in the context only of perfect doubly-nested loops with well-known dependences.

# Going further….

The goal of this segment of the course is to introduce the idea of loop nest transformations, and in particular, the idea of modelling a program's iteration space *geometrically.*

This leads to a large body of powerful and sophisticated techniques based on polyhedral models.  See

- http://en.wikipedia.org/wiki/Polytope_model
- Mohamed-Walid Benabderrahmane, Louis-Noël Pouchet, Albert Cohen, and Cédric Bastoul. 2010. **The polyhedral model is more widely applicable than you think**. In Compiler Construction (CC'10/ETAPS'10), http://dx.doi.org/10.1007/978-3-642-11970-5_16

# Introduction

- In this segment of the course we consider compilation issues for loops involving arrays:

- **How execution order of a loop is constrained,**

- **How a compiler can extract dependence information, and**

- **How this can be used to optimise a program.**

Understanding and transforming execution order can help exploit architectural features:

- **Pipelined, superscalar and VLIW processors**

- **Systems which rely heavily on caches**

- **Processors with special instructions for vectors (SSE, AltiVec)**

- **Multiprocessors, multicore, and co-processors/accelerators**

# Restructuring

- "restructuring compilers" go beyond optimising just to reduce runtime work

- They restructure – in fact *reschedule*

- Conventional optimisations must also be performed
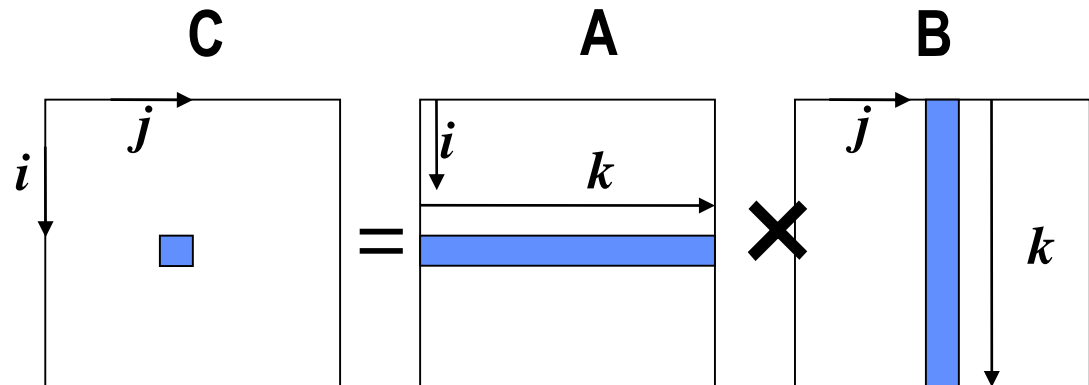- The difference is this:

  - Conventional optimisations reduce the amount of work the computer has to do at run-time

  - Restructuring aims to do the work in an order which suits the target architecture better

# Motivation: an example

```
/*
 * mm: Multiply A by B leaving the
 * result in C.
 * The result matrix is assumed
 * to be initialised to zero.
 */
void mm1(double A[N][N],
         double B[N][N],
         double C[N][N])
{
  int i, j, k;
  for (i = 0; i < N; i++)
    for (j = 0; j < N; j++)
      for (k = 0; k < N; k++)
        C[i][j] += A[i][k] * B[k][j];
}
```
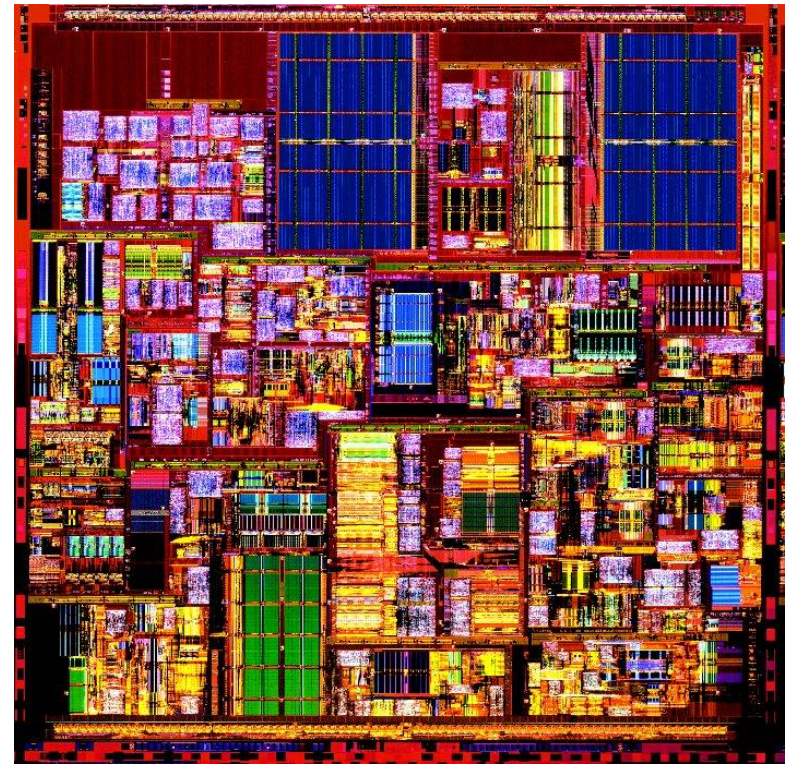
- **We will begin by looking at double-precision floating point matrix multiply**

- **We will investigate the performance of various versions in order to determine what transformations a compiler should apply**
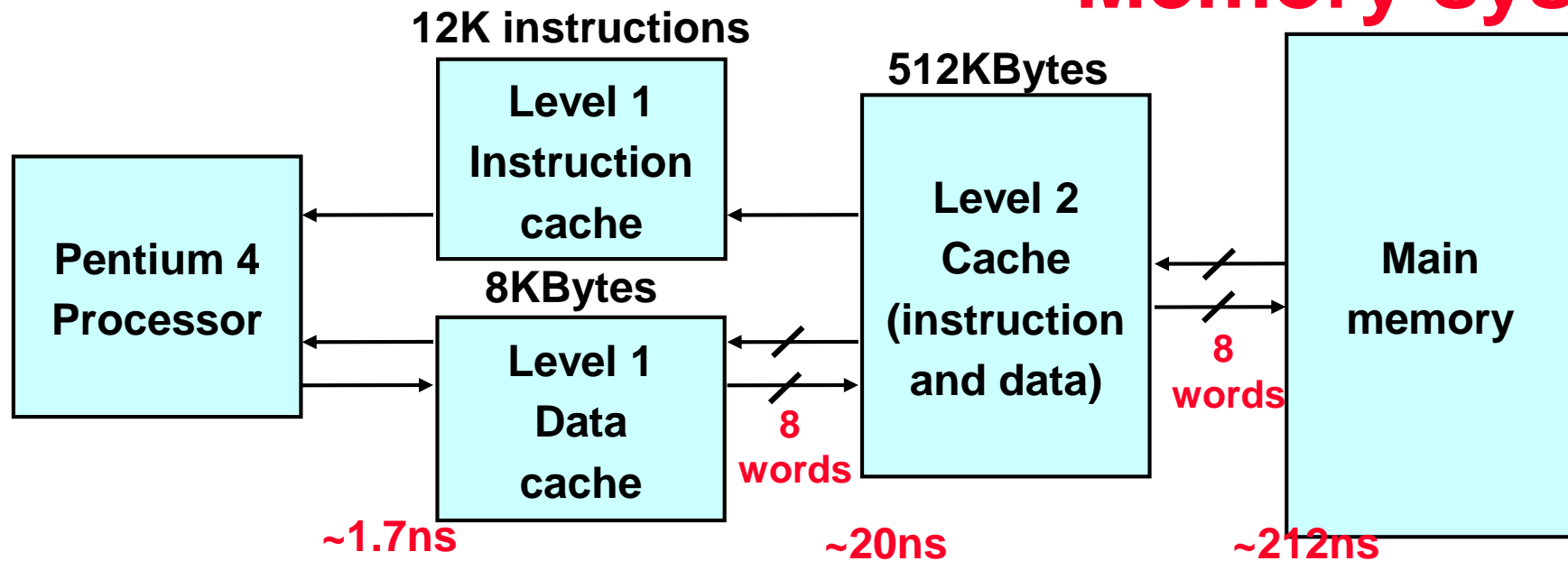
**Outside**



**Inside: Pentium 4 processor**

# Let's experiment with a perfectly ordinary laptop:

- The experiments were performed on a Toshiba Satellite Pro 6100 laptop

- This machine has a 1.6GHz Intel Pentium 4 Mobile processor (we'll look at some other processors shortly)

# Memory system

**12K instructions**

**Level 1 Instruction cache**

**512KBytes**

**Level 2 Cache (instruction and data)**

**Main memory**

**8KBytes**

**Pentium 4 Processor**

**Level 1 Data cache**

**8 words**

**8 words**

**~1.7ns**          **~20ns**          **~212ns**

- L1 Instruction ("trace") cache: 12K microinstructions

- L1 Data cache: 8 KB, 4-way, 64 bytes/line, non-blocking, dual-ported, write-through, pseudo-LRU

- L2 unified cache: 512 KB, 2-Way, 64 Byte/Line, non-blocking

Suppose we're interested in quite big matrices, N=1088

- The matrix occupies $1088^2 \times 8$ = 9.5MBytes

- Each row of the matrix occupies $1088 \times 8$ = 8.5KBytes.

# Performance

- **For N=1088, the initial version runs in 130 seconds.**

- **The matrix multiplication takes $1088^3$ steps, each involving two floating-point operations, an add and a multiply, i.e. $2.6 \times 10^9$ "FLOPs"**

- **This loop achieves a computation rate of 2600/130=19.8 MFLOPs.**

- That is, one floating-point operation completed every 80 clock cycles (the chip runs at 1.6GHz)

- How are we going to get value for money?

# Interchange loops

```
for (i = 0; i < N; i++)
  for (k = 0; k < N; k++)
    for (j = 0; j < N; j++)
    C[i][j] += A[i][k] * B[k][j];
}
```

- **9.6 seconds (267 MFLOPS).**

- **Why is this such a good idea?**

- **How might a compiler perform this transformation?**

- **Does it still give the right output?**

- Can we do better still?

# What was going on?



C        A        B

```
for (i=0; i<N; i++)
  for (j=0; j<N; j++)
    for (k=0; k<N; k++)
      C[i][j]+=A[i][k]*B[k][j]
```

- **IJK variant computes each element of result matrix C one at a time, as inner product of row of A and column of B**
  - **Traverses A in row-major order, B in column-major**



```
for (i=0; i<N; i++)
  for (k=0; k<N; k++)
    for (j=0; j<N; j++)
      C[i][j]+=A[i][k]*B[k][j]
```

- **IKJ variant accumulates partial inner product into a row of result matrix C, using element of A and row of B**
  - **Traverses C and B in row-major order**

# The price of naivety

## Matrix Multiplication Speedup (IKJ/IJK)



- **Relative speedup of IKJ version over IJK version (per machine, per problem size)**

- **On large problems, the IKJ variant is 2-10 times faster**

# Blocking (a.k.a. "tiling")

- **Idea: reorder execution of loop nest so data isn't evicted from cache before it's needed again.**

- **Blocking is a combination of two transformations: "strip mining", followed by interchange; we start with**

```
for (i = 0; i < N; i++)
  for (k = 0; k < N; k++){
    r = A[i][k];
    for (j = 0; j < N; j++)
      C[i][j] += r * B[k][j]; }
```

- **Strip mine the k and j loops:**

```
for (i = 0; i < N; i++)
  for (kk = 0; kk < N; kk += S)
    for (k = kk; k < min(kk+S,N); k++){
      r = A[i][k];
      for (jj = 0; jj < N; jj += S)
        for (j = jj; j < min(jj+S, N); j++)
          C[i][j] += r * B[k][j];
  }
```

# Blocking/tiling – stripmine then interchange

- **Now interchange so blocked loops are outermost:**

```
for (kk = 0; kk < N; kk += S)
  for (jj = 0; jj < N; jj += S)
    for (i = 0; i < N; i++)
      for (k = kk; k < min(kk+S,N); k++){
        r = A[i][k];
        for (j = jj; j < min(jj+S, N); j++)
          C[i][j] += r * B[k][j];
      }
```

- **The inner i,k,j loops perform a multiplication of a pair of partial matrices.**

- **S is chosen so that a S x S submatrix of B and a row of length S of C can fit in the cache.**

- **What is the right value for S?**

# Blocking/tiling – stripmine then interchange

- **Now interchange so blocked loops are outermost:**

```
for (kk = 0; kk < N; kk += S)
  for (jj = 0; jj < N; jj += S)
    for (i = 0; i < N; i++)
      for (k = kk; k < min(kk+S,N); k++){
        r = A[i][k];
        for (j = jj; j < min(jj+S, N); j++)
          C[i][j] += r * B[k][j];
      }
}
```

```
for (kk = 0; kk < N; kk += S)
  for (jj = 0; jj < N; jj += S)
    for (i = 0; i < N; i++)
      C[i][kk:kk+S] +=
        A[i][kk:kk+S]*B[kk:kk+S][jj:jj+S]
```



- **Load a chunk of B into cache; use it to compute the partial inner-products for a column of row-segments of C**

- **Each sweep over row-segment of C multiplies element of A by row segment of submatrix of B.  Sweep again for each A[i][k]**

# Blocking/tiling – stripmine then interchange

**Now interchange so blocked loops are outermost:**

```
for (kk = 0; kk < N; kk += S)
 for (jj = 0; jj < N; jj += S)
   for (i = 0; i < N; i++)
    for (k = kk; k < min(kk+S,N); k++){
     r = A[i][k];
     for (j = jj; j < min(jj+S, N); j++)
      C[i][j] += r * B[k][j];
   }
}
```

```
for (kk = 0; kk < N; kk += S)
 for (jj = 0; jj < N; jj += S)
  for (i = 0; i < N; i++)
   C[i][kk:kk+S] +=
       A[i][kk:kk+S]*B[kk:kk+S][jj:jj+S]
```



**Then we move on to the next block of row segments of C**

No need to revisit the first submatrix of B

# Performance of blocked version: 1.6GHz Pentium 4M (N=1088)



**Problem size 1088**                    **1.6GHz Pentium 4 Mobile (gcc3.4.4)**

**Optimum blocking factor is 144, where we reach 341 MFLOPs**

# Performance of blocked version: Opteron (N=1088)



**Problem size 1088**          **2.4 GHz AMD Opteron (gcc3.4.3) (64KB L1 data cache, 1MB L2)**

**Optimum blocking factor is 64, where we reach 692.4 MFLOPs**

**Since 64 divides 1088 exactly, we can avoid "min" operator, giving 833.6 MFLOPs**

**Using Intel compiler (-Wl,-melf_i386) this reaches 998 MFLOPs**

**Using AMD's AMCL library this machine can reach ~4GFLOP**

# Performance of blocked version: Thinkpad T60 (N=1003)



**1.8 GHz Intel Core2 Duo (Lenovo Thinkpad T60) (gcc3.4.4, 32KB L1 data, 4MB L2)**

**Optimum blocking factor is 48, where we reach 866 MFLOPs**

**Avoiding "min" operator doesn't help.**

**On battery power, clock rate drops to 987MHz, so only 469 MFLOPS (48 still best). In direct proportion to clock rate reduction.**

| Blocking factor | Execution time | MFLOPS |
|---|---|---|
| 8 | 3.815 | 70.4 |
| 16 | 2.784 | 96.4 |
| 32 | 2.283 | 117.6 |
| **40** | **2.193** | **122.4** |
| 48 | 2.253 | 119.1 |
| 56 | 2.473 | 108.5 |
| 64 | 3.404 | 78.9 |
| 72 | 5.608 | 47.9 |
| 80 | 5.578 | 48.1 |
| 88 | 5.808 | 46.2 |
| 96 | 5.928 | 45.3 |
| 104 | 6.309 | 42.5 |
| 112 | 5.778 | 46.5 |



Legend: Execution time/s ■ — MFLOP/s ◆

# Impact….

On Toshiba Satellite Pro 6100 laptop (1.6GHz Pentium 4M):

- Original version: 130 seconds (19.8 MFLOP/s)
- Blocked version: 7.55 seconds (341 MFLOP/s)
  - We started with a "good" optimising compiler!
  - Factor of 17 performance improvement.
  - No reduction in amount of arithmetic performed.


- (Using the Intel library or the ATLAS library does even better)


- See also

- https://github.com/flame/how-to-optimize-gemm/wiki

# Dependence

- Define:
  - IN(S): set of memory locns which might be read by some execn of statement S
  - OUT(S): set of memory locns which might be written by some execn of statement S

- Reordering is constrained by dependences;

- There are four types:
  - Data ("true") dependence: S1 δ S2
    - OUT(S1) ∩ IN(S2)
  - Anti dependence: S1 $\overline{δ}$ S2
    - IN(S1) ∩ OUT(S2)
  - Output dependence: S1 δ° S2
    - OUT(S1) ∩ OUT(S2)
  - Control dependence: S1 δᶜ S2

**("S1 must write something before S2 can read it")**

**("S1 must read something before S2 overwrites it")**

**("If S1 and S2 might both write to a location, S2 must write after S1")**

- These are static analogues of the dynamic RAW, WAR, WAW and control hazards which have to be considered in processor architecture

# Loop-carried dependences

**Consider:**

S1 :  A[0] := 0

      for I = 1 to 8

S2 :  A[I] := A[I-1] + B[I]

**What does this loop do?**

| B: | | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|
| A: | 0 | | | | | | | | |

# Loop-carried dependences

- **Consider:**

S1 : A[0] := 0

for I = 1 to 8

S2 : A[I] := A[I-1] + B[I]

- **What does this loop do?**

| B: | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|

| A: | 0 | 1 | 2 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

- **In this case, there is a data dependence**
  - ➡ **This is a loop-carried dependence - the dependence spans a loop iteration**
  - ➡ **This loop is inherently sequential**

# Loop-carried dependences

- **Consider:**

S1 :    A[0] := 0

   for I = 1 to 8

S2 :    A[I] := A[I-1] + B[I]

- **Loop carried:**

$S2^1$ :    A[1] := A[0] + B[1]

$S2^2$ :    A[2] := A[1] + B[2]

$S2^3$ :    A[3] := A[2] + B[3]

$S2^4$ :    A[4] := A[3] + B[4]

$S2^5$ :    A[5] := A[4] + B[5]

$S2^6$ :    A[6] := A[5] + B[6]

$S2^7$ :    A[7] := A[6] + B[7]

$S2^8$ :    A[8] := A[7] + B[8]

**Dependences cross, from one iteration to next**

0      B[1]
  +
       B[2]
  +
       B[3]
  +
       B[4]
  +
       B[5]
  +
       B[6]
  +
       B[7]
  +
       B[8]
  +

# What is a loop-carried dependence?

- **Consider two iterations $I^1$ and $I^2$**

- **A dependence occurs between two statements $S_p$ and $S_q$ (not necessarily distinct), when an assignment in $S_p{}^{I1}$ refers to the same location as a use in $S_q{}^{I2}$**

➡ **In the example,**

```
S₁ : A[0]  := 0
       for I = 1 to 5
S₂ :    A[I]  := A[I-1] + B[I]
```

- **The assignment is "$A[I^1]$ := ..."**
- **The use is "... := $A[I^2-1]$ ..."**
- **These refer to the same location when $I^1 = I^2-1$**
- **Thus $I^1 < I^2$, ie the assignment is in an earlier iteration**

## Notation: $S_2 \; \delta_< \; S_2$

# **Definition**: The dependence equation

- **A dependence occurs**
  - **between two statements $S_p$ and $S_q$ (not necessarily distinct),**
  - **when there exists a pair of loop iterations $I^1$ and $I^2$,**
  - **such that a memory reference in $S_p$ in $I^1$ may refer to the same location as a memory reference in $S_q$ in $I^2$.**

- **This might occur if $S_p$ and $S_q$ refer to some common array A**

- **Suppose $S_p$ refers to $A[\phi_p(I)]$**       *($\phi_p(I)$ is some subscript*
                                                  *expression involving I)*

- **Suppose $S_q$ refers to $A[\phi_q(I)]$**

- **A dependence of some kind occurs between $S_p$ and $S_q$ if there exists a solution to the equation**

$$\phi_p(I^1) = \phi_q(I^2)$$

- **for integer values of $I^1$ and $I^2$ lying within the loop bounds**

# Types of dependence

- **If a solution to the dependence equation exists, a dependence of some kind occurs**

- **The dependence type depends on what solutions exist**

- The solutions consist of a set of pairs $(I^1, I^2)$

- We would appear to have a *data* dependence if

$$A[\phi_p(I)] \; \epsilon \; OUT(S_p)$$
**and**
$$A[\phi_q(I)] \; \epsilon \; IN(S_q)$$

- But we only really have a data dependence if the assignments *precede* the uses, ie

  - $S_p \; \delta_< \; S_q$
  - if, for each solution pair $(I^1, I^2)$, $I^1 < I^2$

# Dependence versus anti-dependence

- If the *uses* precede the *assignments*, we actually have an *anti*-dependence, ie

$$ S_p \; \overline{\delta}_< \; S_q $$

  if, for each solution pair $(\mathbf{I^1}, \mathbf{I^2})$, $\mathbf{I^1} > \mathbf{I^2}$

- In this case we do have a constraint on execution order

- Because we (may) have to read a value before it (may) be overwritten

- And this anti-dependence is loop-carried

- Anti-dependences prevent re-ordering, and multi-thread parallelism

# Dependence versus anti-dependence

- If there are some solution pairs (I$^1$,I$^2$) with I$^1$ < I$^2$ and some with I$^1$ > I$^2$, we write

$$S_p \ \delta_* \ S_q$$

  This represents that we know we must respect execution ordering, even though the compiler is unable to classify the dependence fully

- If, for all solution pairs (I$^1$,I$^2$), I$^1$ = I$^2$, there are dependences *within* an iteration of the loop, but there are no loop-carried dependences:

$$S_p \ \delta_= \ S_q$$

# Dependence distance

**In many common examples, the set of solution pairs is characterised easily:**

- **Definition:** dependence distance

  - **If, for all solution pairs ($I^1$, $I^2$),**

    $$I^1 = I^2 - k$$

    **then the dependence distance is k**

- **For example in the loop we considered earlier,**

$$S_1 : \texttt{A[0] := 0}$$
$$\texttt{for I = 1 to 5}$$
$$S_2 : \texttt{A[I] := A[I-1] + B[I]}$$

**We find that $S_2\ \delta_<\ S_2$ with dependence distance 1.**

- **_((of course there are many cases where the difference is not constant and so the dependence cannot be summarised this way))._**

# Reuse distance

- When optimising for cache performance, it is sometimes useful to consider the re-use relationship,

  - $IN(S_1) \cap IN(S_2)$

- Here there is no dependence - it doesn't matter which read occurs first
- Nonetheless, cache performance can be improved by minimising the *reuse distance*

- The reuse distance is calculated essentially the same way

- Eg

  for I = 5 to 100

  S1:  B[I] := A[I] * 2

  S2:  C[I] := A[I-5] * 10

- Here we have a loop-carried reuse with distance 5

- Up to now we have looked at single loops

- Now let's generalise to loop "nests"

- We begin by considering a very common dependence pattern, called the "wavefront":

$$\texttt{for } I_1 \texttt{ = 0 to 3 do}$$
$$\quad \texttt{for } I_2 \texttt{ = 0 to 3 do}$$
$$S: \quad \texttt{A[}I_1,I_2\texttt{] := A[}I_1-1,I_2\texttt{] + A[}I_1,I_2-1\texttt{]}$$

- Dependence structure?

# System of dependence equations

**Consider the dependence equations for this loop nest:**

```
for I₁ = 0 to 3 do
   for I₂ = 0 to 3 do
```
$$S: \quad \text{A}[I_1, I_2] \;:=\; \text{A}[I_1 - 1, I_2] \;+\; \text{A}[I_1, I_2 - 1]$$

**There are two potential dependences arising from the three references to A, so two systems of dependence equations to solve:**

1. Between $\text{A}[I_1^1, I_2^1]$ and $\text{A}[I_1^2 - 1, I_2^2]$:

$$\begin{cases} I_1^1 &= I_1^2 - 1 \\ I_2^1 &= I_2^2 \end{cases}$$

2. Between $\text{A}[I_1^1, I_2^1]$ and $\text{A}[I_1^2, I_2^2 - 1]$:

$$\begin{cases} I_1^1 = I_1^2 \\ I_2^1 = I_2^2 - 1 \end{cases}$$

- **The same loop:**

```
for I₁ = 0 to 3 do
    for I₂ = 0 to 3 do
```
$$S: \quad \text{A}[I_1,I_2] := \text{A}[I_1-1,I_2] + \text{A}[I_1,I_2-1]$$

- **For humans the easy way to understand this loop nest is to draw the *iteration space graph* showing the iteration-to-iteration dependences:**

$$
\begin{array}{cccccccc}
S^{00} & \xrightarrow{\delta} & S^{01} & \xrightarrow{\delta} & S^{02} & \xrightarrow{\delta} & S^{03} \\
\delta \downarrow & & \delta \downarrow & & \delta \downarrow & & \delta \downarrow \\
S^{10} & \xrightarrow{\delta} & S^{11} & \xrightarrow{\delta} & S^{12} & \xrightarrow{\delta} & S^{13} \\
\delta \downarrow & & \delta \downarrow & & \delta \downarrow & & \delta \downarrow \\
S^{20} & \xrightarrow{\delta} & S^{21} & \xrightarrow{\delta} & S^{22} & \xrightarrow{\delta} & S^{23} \\
\delta \downarrow & & \delta \downarrow & & \delta \downarrow & & \delta \downarrow \\
S^{30} & \xrightarrow{\delta} & S^{31} & \xrightarrow{\delta} & S^{32} & \xrightarrow{\delta} & S^{33}
\end{array}
$$

- **The diagram shows an arrow for each solution of each dependence equation.  Is there any parallelism?**

$$S^{00} \xrightarrow{\delta} S^{01} \xrightarrow{\delta} S^{02} \xrightarrow{\delta} S^{03}$$

$$\delta \downarrow \quad\quad \delta \downarrow \quad\quad \delta \downarrow \quad\quad \delta \downarrow$$

$$S^{10} \xrightarrow{\delta} S^{11} \xrightarrow{\delta} S^{12} \xrightarrow{\delta} S^{13}$$

$$\delta \downarrow \quad\quad \delta \downarrow \quad\quad \delta \downarrow \quad\quad \delta \downarrow$$

$$S^{20} \xrightarrow{\delta} S^{21} \xrightarrow{\delta} S^{22} \xrightarrow{\delta} S^{23}$$

$$\delta \downarrow \quad\quad \delta \downarrow \quad\quad \delta \downarrow \quad\quad \delta \downarrow$$

$$S^{30} \xrightarrow{\delta} S^{31} \xrightarrow{\delta} S^{32} \xrightarrow{\delta} S^{33}$$

```
for I₁ = 0 to 3 do
  for I₂ = 0 to 3 do
S:   A[I₁,I₂] := A[I₁ - 1,I₂] + A[I₁,I₂ - 1]
```

- **The inner loop is not vectorisable since there is a dependence chain linking successive iterations.**

  - **(to use a vector instruction, need to be able to operate on each element of the vector in parallel)**

- **Similarly, the outer loop is not parallel**

- **This loop is *interchangeable*: the top-to-bottom, left-to-right execution order is also valid since all dependence constraints (as shown by the arrows) are still satisfied.**

- **Interchanging the loop does not improve vectorisability or parallelisability**

$S^{00}$ $\xrightarrow{\delta}$ $S^{01}$ $\xrightarrow{\delta}$ $S^{02}$ $\xrightarrow{\delta}$ $S^{03}$

$\delta \downarrow$     $\delta \downarrow$     $\delta \downarrow$     $\delta \downarrow$

$S^{10}$ $\xrightarrow{\delta}$ $S^{11}$ $\xrightarrow{\delta}$ $S^{12}$ $\xrightarrow{\delta}$ $S^{13}$

$\delta \downarrow$     $\delta \downarrow$     $\delta \downarrow$     $\delta \downarrow$

$S^{20}$ $\xrightarrow{\delta}$ $S^{21}$ $\xrightarrow{\delta}$ $S^{22}$ $\xrightarrow{\delta}$ $S^{23}$

$\delta \downarrow$     $\delta \downarrow$     $\delta \downarrow$     $\delta \downarrow$

$S^{30}$ $\xrightarrow{\delta}$ $S^{31}$ $\xrightarrow{\delta}$ $S^{32}$ $\xrightarrow{\delta}$ $S^{33}$

```
for I_1 = 0 to 3 do
  for I_2 = 0 to 3 do
S:   A[I_1,I_2] := A[I_1 - 1,I_2] + A[I_1,I_2 - 1]
```

$$\text{for } I_1 = 0 \text{ to } 3 \text{ do}$$
$$\text{for } I_2 = 0 \text{ to } 3 \text{ do}$$
$$S: \quad A[I_1,I_2] := A[I_1 - 1,I_2] + A[I_1,I_2 - 1]$$

- **The inner loop is not vectorisable since there is a dependence chain linking successive iterations.**

  - **(to use a vector instruction, need to be able to operate on each element of the vector in parallel)**

- **Similarly, the outer loop is not parallel**

- **This loop is *interchangeable*: the top-to-bottom, left-to-right execution order is also valid since all dependence constraints (as shown by the arrows) are still satisfied.**

- **Interchanging the loop does not improve vectorisability or parallelisability**

$$S^{00} \xrightarrow{\delta} S^{01} \xrightarrow{\delta} S^{02} \xrightarrow{\delta} S^{03}$$
$$\delta\downarrow \qquad \delta\downarrow \qquad \delta\downarrow \qquad \delta\downarrow$$
$$S^{10} \xrightarrow{\delta} S^{11} \xrightarrow{\delta} S^{12} \xrightarrow{\delta} S^{13}$$
$$\delta\downarrow \qquad \delta\downarrow \qquad \delta\downarrow \qquad \delta\downarrow$$
$$S^{20} \xrightarrow{\delta} S^{21} \xrightarrow{\delta} S^{22} \xrightarrow{\delta} S^{23}$$
$$\delta\downarrow \qquad \delta\downarrow \qquad \delta\downarrow \qquad \delta\downarrow$$
$$S^{30} \xrightarrow{\delta} S^{31} \xrightarrow{\delta} S^{32} \xrightarrow{\delta} S^{33}$$
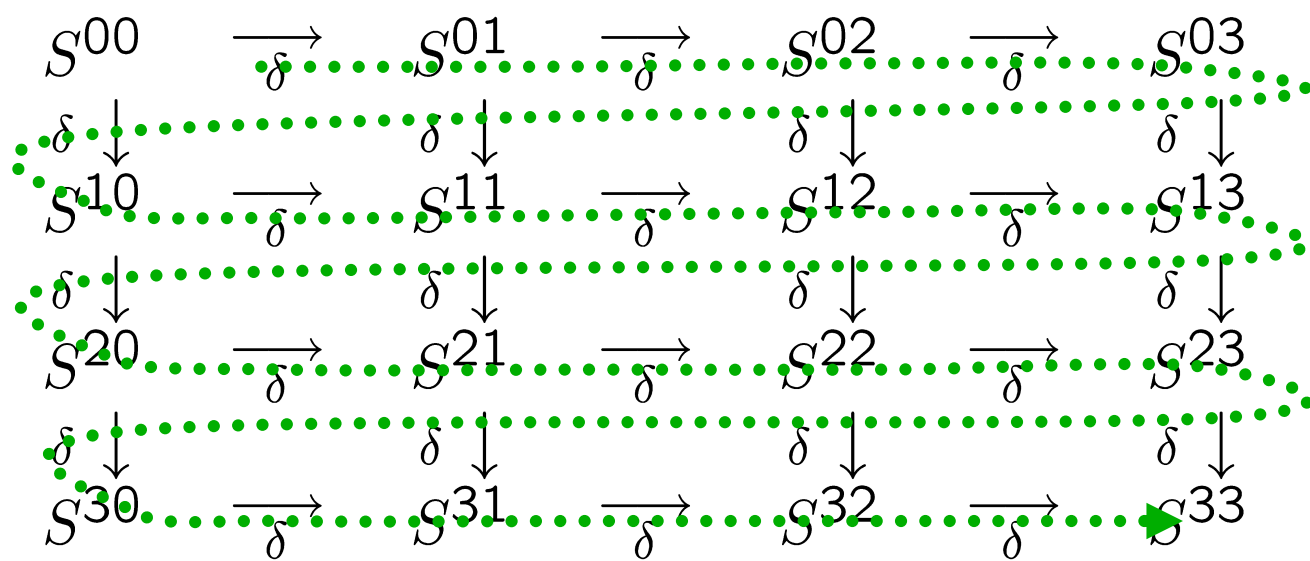
```
for I_1 = 0 to 3 do
  for I_2 = 0 to 3 do
S:   A[I_1,I_2] := A[I_1 - 1,I_2] + A[I_1,I_2 - 1]
```

- ▶ **The inner loop is not vectorisable since there is a dependence chain linking successive iterations.**
  - ➡ **(to use a vector instruction, need to be able to operate on each element of the vector in parallel)**

- · **Similarly, the outer loop is not parallel**

- · **This loop is *interchangeable*: the top-to-bottom, left-to-right execution order is also valid since all dependence constraints (as shown by the arrows) are still satisfied.**

- · **Interchanging the loop does not improve vectorisability or parallelisability**

# Interchange: counter-example

```
for I₁ = 0 to 3 do
  for I₂ = 0 to 3 do
    A[I₁,I₂] := A[I₁ - 1, I₂ + 1] + B[I₁, I₂]
```

```
for I₁ = 0 to 3 do
  for I₂ = 0 to 3 do
    A[I₁,I₂] := A[I₁ - 1, I₂ + 1] + B[I₁, I₂]
```

$$S^{00} \quad\quad S^{01} \quad\quad S^{02} \quad\quad S^{03}$$

$$S^{10} \quad\quad S^{11} \quad\quad S^{12} \quad\quad S^{13}$$

$$S^{20} \quad\quad S^{21} \quad\quad S^{22} \quad\quad S^{23}$$

$$S^{30} \quad\quad S^{31} \quad\quad S^{32} \quad\quad S^{33}$$

# Interchange: counter-example

```
for I₁ = 0 to 3 do
   for I₂ = 0 to 3 do
      A[I₁,I₂] := A[I₁ - 1, I₂ + 1] + B[I₁, I₂]
```

**Before interchange**

$$S^{00} \qquad S^{01} \qquad S^{02} \qquad S^{03}$$

$$S^{10} \qquad S^{11} \qquad S^{12} \qquad S^{13}$$

$$S^{20} \qquad S^{21} \qquad S^{22} \qquad S^{23}$$

$$S^{30} \qquad S^{31} \qquad S^{32} \qquad S^{33}$$

# Interchange: counter-example

```
for I_1 = 0 to 3 do
  for I_2 = 0 to 3 do
    A[I_1,I_2] := A[I_1 - 1, I_2 + 1] + B[I_1, I_2]
```

**After interchange:**

**New traversal order crosses dependence arrows backwards**

$$S^{00} \quad S^{01} \quad S^{02} \quad S^{03}$$

$$S^{10} \quad S^{11} \quad S^{12} \quad S^{13}$$

$$S^{20} \quad S^{21} \quad S^{22} \quad S^{23}$$

$$S^{30} \quad S^{31} \quad S^{32} \quad S^{33}$$

# Interchange: condition

- A loop is *interchangeable* if all dependence constraints (as shown by the arrows) are still satisfied by the top-to-bottom, left-to-right execution order

- How can you tell whether a loop can be interchanged?

- Look at its dependence direction vectors:
  - Is there a dependence direction vector with the form (<,>) ?
    - ie there is a dependence distance vector $(k_1, k_2)$ with $k_1 > 0$ and $k_2 < 0$ ?
    - If so, interchange would be *invalid*
  - Because the arrows would be traversed backwards
  - All other dependence directions are OK.

**Consider this variation on the wavefront loop:**

```
for k₁ := 0 to 3 do
   for k₂ := k₁ to k₁+3 do
```
$S$ : A[k₁,k₂-k₁] := A[k₁-1,k₂-k₁]+A[k₁,k₂-k₁-1]

- The inner loop's control variable runs from $k_1$ to $k_1+3$.

- The iteration space of this loop has $4^2$ iterations just like the original loop.

- If we draw the iteration space with each iteration $S^{K_1,K_2}$ at coordinate position $(K_1,K_2)$, it is skewed to form a lozenge shape:

$$S^{00} \quad S^{01} \quad S^{02} \quad S^{03}$$
$$S^{11} \quad S^{12} \quad S^{13} \quad S^{14}$$
$$S^{22} \quad S^{23} \quad S^{24} \quad S^{25}$$
$$S^{33} \quad S^{34} \quad S^{35} \quad S^{36}$$

**This loop performs the same computation as the original.**

- Thus the dependence structure of the skewed loop is shown by marking the iteration space with all the dependences:

$$S^{00} \xrightarrow{\delta} S^{01} \xrightarrow{\delta} S^{02} \xrightarrow{\delta} S^{03}$$

$$\searrow{\delta} \quad \searrow{\delta} \quad \searrow{\delta} \quad \searrow{\delta}$$

$$S^{11} \xrightarrow{\delta} S^{12} \xrightarrow{\delta} S^{13} \xrightarrow{\delta} S^{14}$$

$$\searrow{\delta} \quad \searrow{\delta} \quad \searrow{\delta} \quad \searrow{\delta}$$

$$S^{22} \xrightarrow{\delta} S^{23} \xrightarrow{\delta} S^{24} \xrightarrow{\delta} S^{25}$$

$$\searrow{\delta} \quad \searrow{\delta} \quad \searrow{\delta} \quad \searrow{\delta}$$

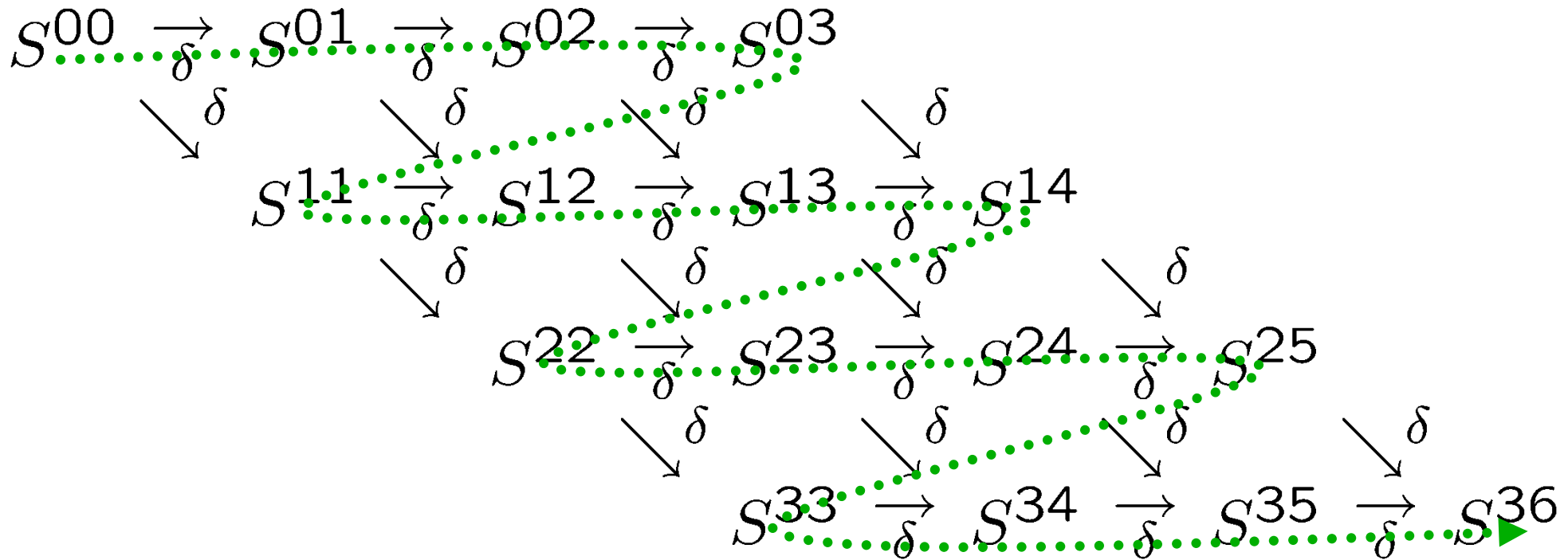$$S^{33} \xrightarrow{\delta} S^{34} \xrightarrow{\delta} S^{35} \xrightarrow{\delta} S^{36}$$

- Can this loop nest be vectorised?
- Can this loop nest be interchanged?

# Skewing changes effect of interchange

- Thus the dependence structure of the skewed loop is shown by marking the iteration space with all the dependences:

$$S^{00} \xrightarrow{\delta} S^{01} \xrightarrow{\delta} S^{02} \xrightarrow{\delta} S^{03}$$

$$S^{11} \xrightarrow{\delta} S^{12} \xrightarrow{\delta} S^{13} \xrightarrow{\delta} S^{14}$$

$$S^{22} \xrightarrow{\delta} S^{23} \xrightarrow{\delta} S^{24} \xrightarrow{\delta} S^{25}$$

$$S^{33} \xrightarrow{\delta} S^{34} \xrightarrow{\delta} S^{35} \xrightarrow{\delta} S^{36}$$
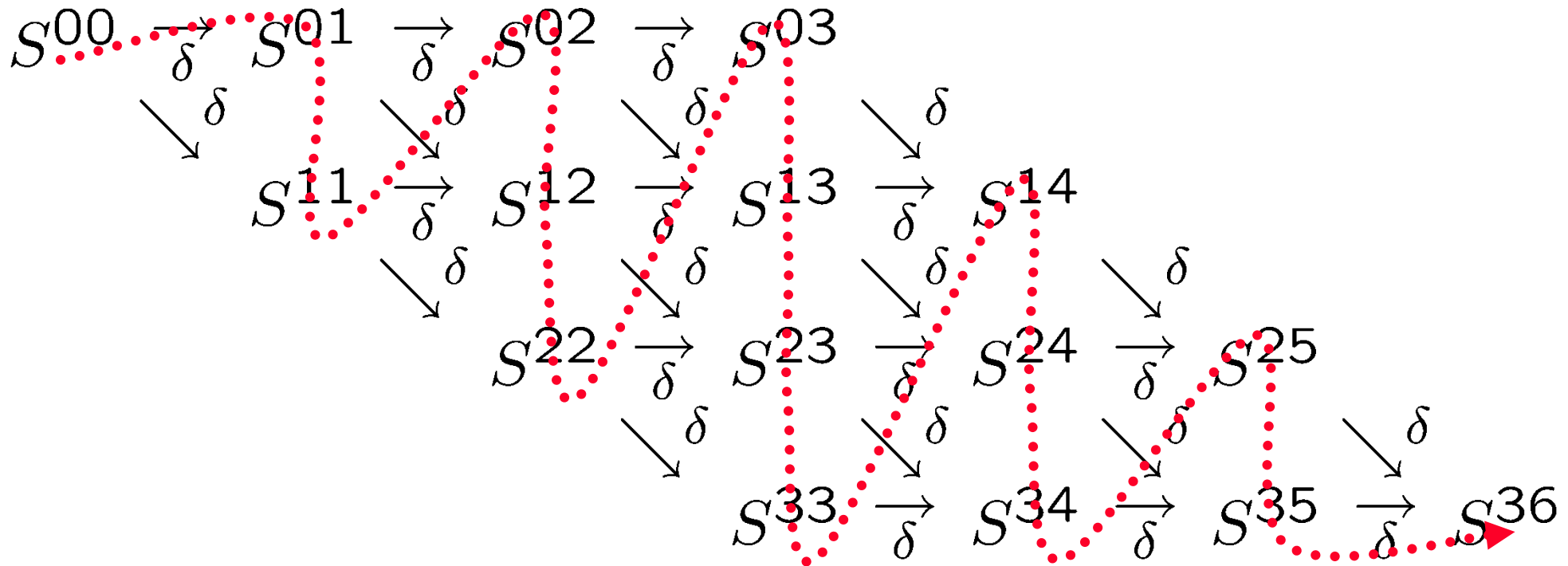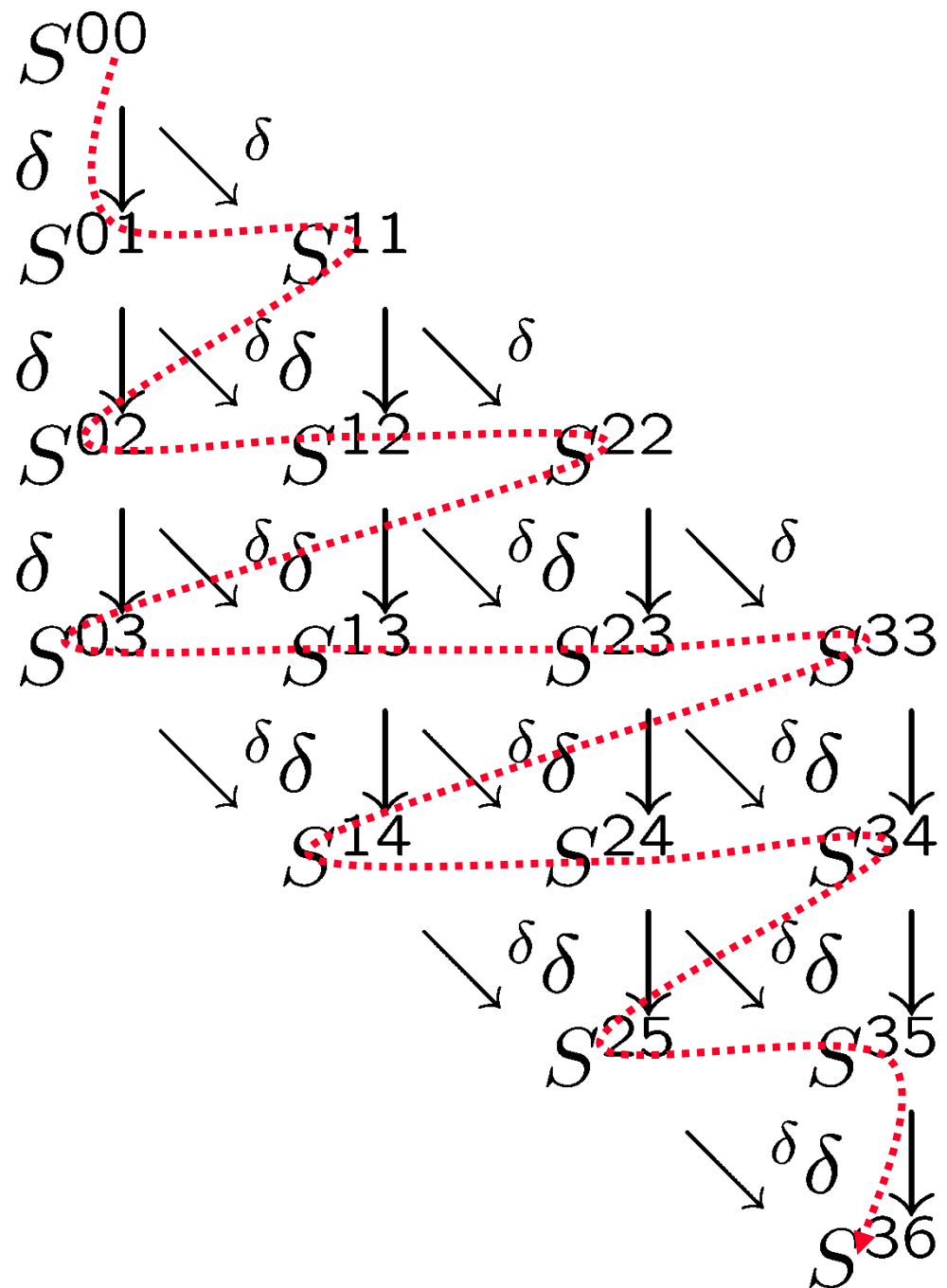
- Original execution order

# Interchange after skewing

- Thus the dependence structure of the skewed loop is shown by marking the iteration space with all the dependences:

$$S^{00} \xrightarrow{\delta} S^{01} \xrightarrow{\delta} S^{02} \xrightarrow{\delta} S^{03}$$

$$\searrow \delta \qquad \searrow \delta \qquad \searrow \delta \qquad \searrow \delta$$

$$S^{11} \xrightarrow{\delta} S^{12} \xrightarrow{\delta} S^{13} \xrightarrow{\delta} S^{14}$$

$$\searrow \delta \qquad \searrow \delta \qquad \searrow \delta \qquad \searrow \delta$$

$$S^{22} \xrightarrow{\delta} S^{23} \xrightarrow{\delta} S^{24} \xrightarrow{\delta} S^{25}$$

$$\searrow \delta \qquad \searrow \delta \qquad \searrow \delta \qquad \searrow \delta$$

$$S^{33} \xrightarrow{\delta} S^{34} \xrightarrow{\delta} S^{35} \xrightarrow{\delta} S^{36}$$

- Transposed execution order

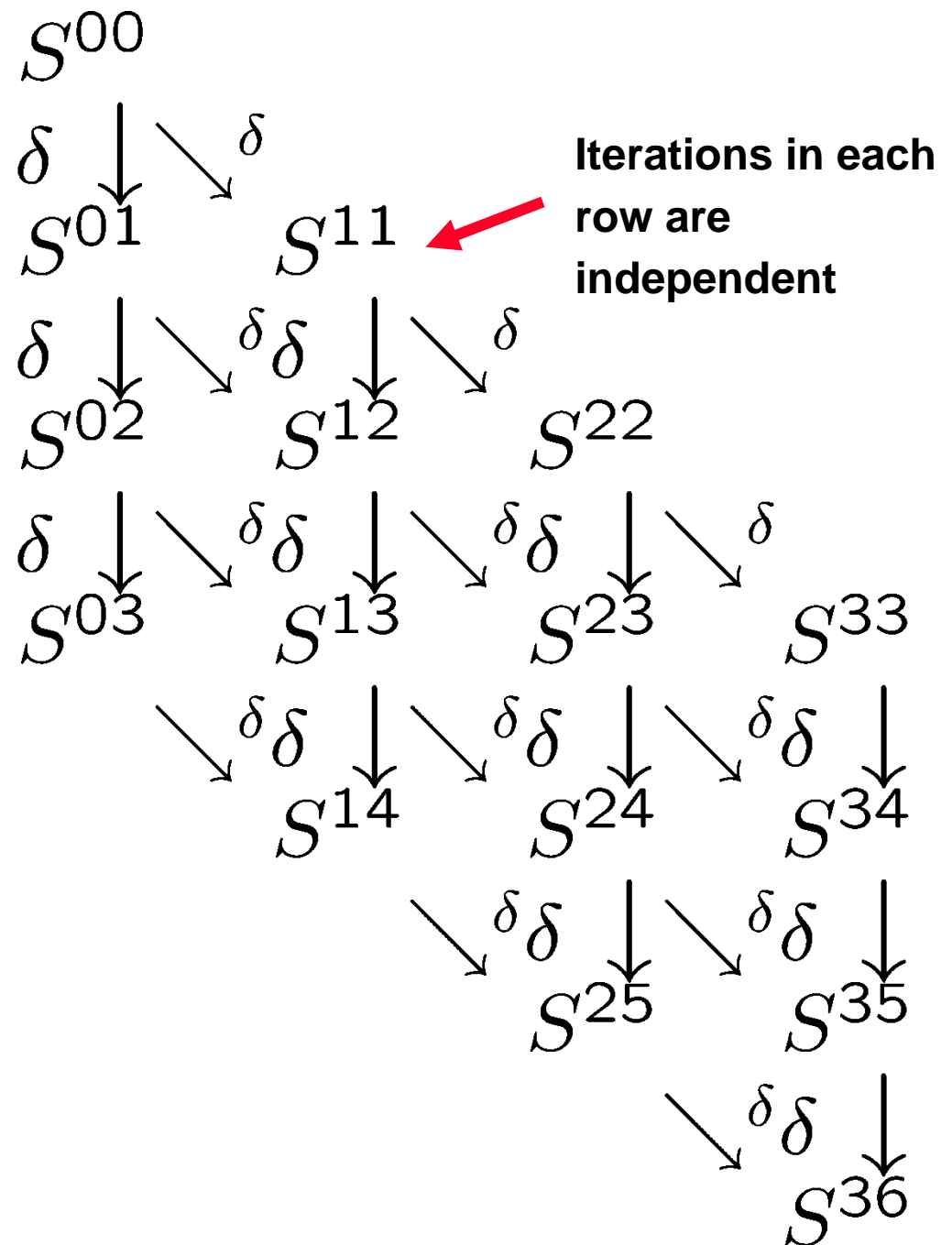- You can think of loop interchange as changing the way the iteration space is traversed

- Alternatively, you can think of it as a change to the way the runtime code instances are mapped onto the iteration space

- Traversal is always lexicographic – ie left-to-right, top-down

$S^{00}$

$\delta$ $\delta$

$S^{01}$ $S^{11}$

$\delta$ $\delta$ $\delta$ $\delta$

$S^{02}$ $S^{12}$ $S^{22}$

$\delta$ $\delta$ $\delta$ $\delta$ $\delta$ $\delta$

$S^{03}$ $S^{13}$ $S^{23}$ $S^{33}$

$\delta$ $\delta$ $\delta$ $\delta$ $\delta$ $\delta$

$S^{14}$ $S^{24}$ $S^{34}$

$\delta$ $\delta$ $\delta$ $\delta$
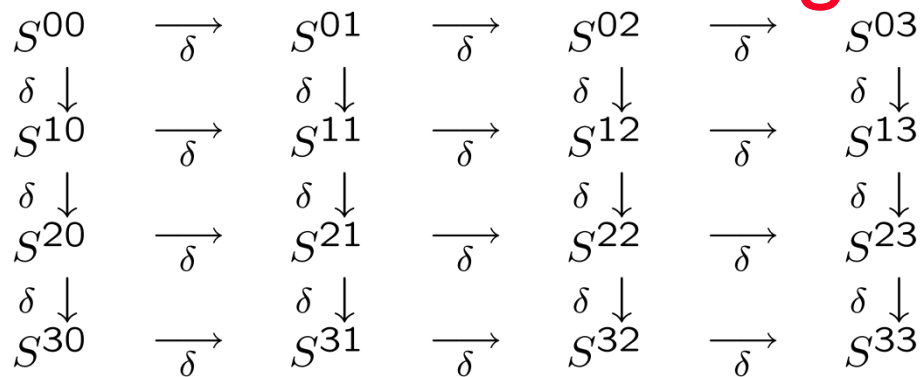
$S^{25}$ $S^{35}$

$\delta$ $\delta$

$S^{36}$

- The inner loop is now vectorisable, since it has no loop-carried dependence

- The skewed iteration space has N rows and 2N-1 columns, but still only $N^2$ actual statement instances.

$S^{00}$

$\delta \downarrow \quad \searrow \delta$

$S^{01} \qquad S^{11}$ ← **Iterations in each row are independent**

$\delta \downarrow \searrow \delta \delta \downarrow \searrow \delta$

$S^{02} \qquad S^{12} \qquad S^{22}$

$\delta \downarrow \searrow \delta \delta \downarrow \searrow \delta \delta \downarrow \searrow \delta$

$S^{03} \qquad S^{13} \qquad S^{23} \qquad S^{33}$

$\searrow \delta \delta \downarrow \searrow \delta \delta \downarrow \searrow \delta \delta \downarrow$

$S^{14} \qquad S^{24} \qquad S^{34}$

$\searrow \delta \delta \downarrow \searrow \delta \delta \downarrow$

$S^{25} \qquad S^{35}$

$\searrow \delta \delta \downarrow$

$S^{36}$

$$S^{00} \xrightarrow{\delta} S^{01} \xrightarrow{\delta} S^{02} \xrightarrow{\delta} S^{03}$$
$$\delta \downarrow \qquad \delta \downarrow \qquad \delta \downarrow \qquad \delta \downarrow$$
$$S^{10} \xrightarrow{\delta} S^{11} \xrightarrow{\delta} S^{12} \xrightarrow{\delta} S^{13}$$
$$\delta \downarrow \qquad \delta \downarrow \qquad \delta \downarrow \qquad \delta \downarrow$$
$$S^{20} \xrightarrow{\delta} S^{21} \xrightarrow{\delta} S^{22} \xrightarrow{\delta} S^{23}$$
$$\delta \downarrow \qquad \delta \downarrow \qquad \delta \downarrow \qquad \delta \downarrow$$
$$S^{30} \xrightarrow{\delta} S^{31} \xrightarrow{\delta} S^{32} \xrightarrow{\delta} S^{33}$$

```
for I_1 = 0 to 3 do
   for I_2 = 0 to 3 do
S :    A[I_1,I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]
```

- **Original loop interchangeable but not vectorisable.**
- **We skewed inner loop by outer loop by factor 1.**
- **Still not vectorisable, but interchangeable.**
- **Interchanged, skewed loop *is* vectorisable.**
- **Bounds of new loop not simple!**

$$S^{00}$$
$$\delta \downarrow \searrow \delta$$
$$S^{01} \qquad S^{11}$$
$$\delta \downarrow \searrow \delta \; \delta \downarrow \searrow \delta$$
$$S^{02} \qquad S^{12} \qquad S^{22}$$
$$\delta \downarrow \searrow \delta \; \delta \downarrow \searrow \delta \; \delta \downarrow \searrow \delta$$
$$S^{03} \qquad S^{13} \qquad S^{23} \qquad S^{33}$$
$$\searrow \delta \; \delta \downarrow \searrow \delta \; \delta \downarrow \searrow \delta \; \delta \downarrow$$
$$S^{14} \qquad S^{24} \qquad S^{34}$$
$$\searrow \delta \; \delta \downarrow \searrow \delta \; \delta \downarrow$$
$$S^{25} \qquad S^{35}$$
$$\searrow \delta \; \delta \downarrow$$
$$S^{36}$$

```
for k_2 := 0 to 2N_2 - 2 do
   for k_1 := max(0, K_2 - N_2 + 2) to min(K_2, N_1) do
S : A[k_1,k_2-k_1] := A[k_1-1,k_2-k_1]+A[k_1,k_2-k_1-1]
```

- **Is skewing ever invalid?**
- **Does skewing affect interchangeability?**
- **Does skewing affect dependence distances?**
- **Can you predict value of skewing?**

# Summary: dependence

- Dependence equation for single loop:

  - **Suppose $S_p$ refers to $A[\phi_p(I)]$**

  - **Suppose $S_q$ refers to $A[\phi_q(I)]$**

  - **A dependence of some kind occurs between $S_p$ and $S_q$ if there exists a solution to the equation**

  $$\phi_p(I^1) = \phi_q(I^2)$$

  - **for integer values of $I^1$ and $I^2$ lying within the loop bounds**

- For doubly-nested loops **over multidimensional arrays, generalise to system of simultaneous dependence equations for two iterations, $(I_1{}^1, I_2{}^1)$ and $(I_1{}^2, I_2{}^2)$**

- Iteration space graph, **lexicographic schedule of execution**

- **Arrows in graph show solutions to dependence equation**

- Dependence distance vectors **characterise families of congruent arrows**

# Summary: transformations

- A loop can be executed in parallel **if it has no loop-carried dependence**

- A loop nest can be interchanged **if the transposed dependence distance vectors are lexicographically forward**

- Strip-mining **is always valid**

- Tiling **= strip-mining + interchange**

- **Skewing is always valid**

- Skewing can expose parallelism **by aligning parallel iterations with one of the loops**

- **Skewing can make interchange (and therefore tiling) valid**

# Matrix representation of loop transformations

- **To skew the inner loop by the outer loop by factor 1 we adjust the loop bounds, and replace $I_1$ by $K_1$, and $I_2$ by $K_2$-$K_1$. That is,**

    **$(K_1,K_2) = (I_1,I_2) \cdot U$**

- **where U is a 2 x 2 matrix**

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

- **That is,**

    **$(K_1,K_2) = (I_1,I_2) \cdot U = (I_1,I_2+I_1)$**

- **The inverse gets us back again:**

    **$(I_1,I_2) = (K_1,K_2) \cdot U^{-1} = (K_1,K_2-K_1)$**

- Matrix U maps each statement instance $S^{I_1 I_2}$ to its position in the new iteration space, $S^{K_1 K_2}$:

- Original iteration space:

$I_2$ :

|       | 0        | 1        | 2        | 3        |
|-------|----------|----------|----------|----------|
| $I_1$ |          |          |          |          |
| 0     | $S^{00}$ | $S^{01}$ | $S^{02}$ | $S^{03}$ |
| 1     | $S^{10}$ | $S^{11}$ | $S^{12}$ | $S^{13}$ |
| 2     | $S^{20}$ | $S^{21}$ | $S^{22}$ | $S^{23}$ |
| 3     | $S^{30}$ | $S^{31}$ | $S^{32}$ | $S^{33}$ |

```
for I_1 = 0 to 3 do
    for I_2 = 0 to 3 do
S:     A[I_1,I_2] := A[I_1-1,I_2] + A[I_1,I_2-1]
```

⬇

```
for k_1 := 0 to 3 do
    for k_2 := k_1 to k_1+3 do
S: A[k_1,k_2-k_1] := A[k_1-1,k_2-k_1]+A[k_1,k_2-k_1-1]
```

**The subscripts are mapped back using U⁻¹**

$K_2$ :

|       | 0        | 1        | 2        | 3        | 4        | 5        | 6        |
|-------|----------|----------|----------|----------|----------|----------|----------|
| $K_1$ |          |          |          |          |          |          |          |
| 0     | $S^{00}$ | $S^{01}$ | $S^{02}$ | $S^{03}$ |          |          |          |
| 1     |          | $S^{11}$ | $S^{12}$ | $S^{13}$ | $S^{14}$ |          |          |
| 2     |          |          | $S^{22}$ | $S^{23}$ | $S^{24}$ | $S^{25}$ |          |
| 3     |          |          |          | $S^{33}$ | $S^{34}$ | $S^{35}$ | $S^{36}$ |

**The dependences are subject to the same transformation.**

$(K_1,K_2) = (I_1,I_2) \cdot U = (I_1,I_2+I_1)$

$(I_1,I_2) = (K_1,K_2) \cdot U^{-1} = (K_1,K_2-K_1)$

# Using matrices to reason about dependence

Recall that:

- There is a dependence between two iterations $(I_1^1, I_2^1)$ and $(I_1^2, I_2^2)$ if there is a memory location which is assigned to in iteration $(I_1^1, I_2^1)$, and read in iteration $(I_1^2, I_2^2)$.

  ((unless there is an intervening assignment))

- If $(I_1^1, I_2^1)$ precedes $(I_1^2, I_2^2)$ it is a *data*-dependence.
- If $(I_1^2, I_2^2)$ precedes $(I_1^1, I_2^1)$ it is a *anti*-dependence.
- If the location is assigned to in both iterations, it is an *output*-dependence.

- The dependence distance vector $(D_1, D_2)$ is $(I_1^1 - I_1^2, I_2^1 - I_2^2)$.

# Transforming dependence vectors

- If there is a dependence between two iterations $(I_1^1, I_2^1)$ and $(I_1^2, I_2^2)$

- Then iterations $(I_1^1, I_2^1) \cdot U$ and $(I_1^2, I_2^2) \cdot U$ will also read and write the same location

- The transformation U is *valid* iff

  $(I_1^1, I_2^1) \cdot U$ precedes $(I_1^2, I_2^2) \cdot U$

  whenever there is a dependence between

  $(I_1^1, I_2^1)$ and $(I_1^2, I_2^2)$.

- In the transformed loop the dependence distance vector is also transformed, to

  $(D_1, D_2) \cdot U$

# Transforming dependence vectors

- **What do we mean by "precedes"?**

**Definition: Lexicographic ordering:**

$(I_1^1, I_2^1)$ **precedes** $(I_1^2, I_2^2)$

**if** $I_1^1 < I_1^2$, **or** $I_1^1 = I_1^2$ **and** $I_2^1 < I_2^2$

- **"Lexicographic" is dictionary order – both "baz" and "can" precede "cat"**

- **So (1,2) precedes (1,3)**

- **But (0,3) precedes (1,4)**

Before transformation we had two dependences:

1. Distance: (1,0), direction: (<,.)

2. Distance: (0,1), direction: (.,<)

- After transformation by matrix

$$U = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

- (i.e. skewing of inner loop by outer) we get:

1. Distance: (1,1), direction: (<,<)

2. Distance: (0,1), direction: (.,<)

- We can also represent loop interchange by a matrix transformation.
- After transforming the skewed loop by matrix

$$V = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

- (i.e. loop interchange) we get:

1. Distance: (1,1), direction: (<,<)

2. Distance: (1,0), direction: (<,.)

- The transformed iteration space is the transpose of the skewed iteration space:

$$
\begin{array}{cccc}
S^{00} & & & \\
S^{10} & S^{11} & & \\
S^{20} & S^{21} & S^{22} & \\
S^{30} & S^{31} & S^{32} & S^{33} \\
& S^{41} & S^{42} & S^{43} \\
& & S^{52} & S^{53} \\
& & & S^{63}
\end{array}
$$

# Summary

- $(I_1,I_2) \cdot U$ maps each statement instance $(I_1,I_2)$ to its new position $(K_1,K_2)$ in the transformed loop's execution sequence

- $(D_1,D_2) \cdot U$ gives new dependence distance vector, giving test for validity

- Captures skewing, interchange and reversal

- Compose transformations by matrix multiplication

    $U_1 \cdot U_2$

- Resulting loop's bounds may be a little tricky

    - Efficient algorithms exist [Banerjee90] to maximise parallelism by skewing and loop interchanging

    - Efficient algorithms exist to optimise cache performance by finding the combination of blocking, block size, interchange and skewing which leads to the best reuse [Wolf91]

Restructuring compilers - conclusions:

- Restructuring compilers can find parallelism

- And enhance locality

- **For a very restricted class of programs**

  - For-loops over arrays with array subscripts that are simple ("affine") expressions involving loop control variables

- But for this restricted class there is a rather elegant theory (the "polyhedral" or "polytope" model, http://en.wikipedia.org/wiki/Polytope_model)

- Extending beyond this is a big research problem

- Current compilers (GCC, Intel, Microsoft etc) can do some of this, in theory – but are often defeated by program complexity

# References

- **Hennessy and Patterson: Section 4.4 (pp.319)**
- **Background: "conventional" compiler techniques**
  - A.V. Aho, R. Sethi, and J.D. Ullman. Compilers: Principles, Techniques and Tools. Addison Wesley, 1986.
  - Andrew Appel and Jens Palsberg, Modern Compiler Implementation. Cambridge University Press, 2002.
  - Cooper and Torczon, Engineering a Compiler. Morgan Kaufmann 2004.
  - Morgan, Building an Optimizing Compiler
- **Textbooks covering restructuring compilers**
  - Michael Wolfe. High Performance Compilers for Parallel Computing. Addison Wesley, 1996.
  - Steven Muchnick, Advanced Compiler Design and Implementation. Morgan Kaufmann, 1997.
  - Ken Kennedy and Randy Allen, Optimizing Compilers for Modern Architectures. Morgan Kaufmann, 2001.
- **Research papers:**

  **READ THIS ONE**
  - D. F. Bacon and S. L. Graham and O. J. Sharp, "Compiler Transformations for High-Performance Computing". ACM Computing Surveys V26 N4 Dec 1994
    http://doi.acm.org/10.1145/197405.197406
  - U. Banerjee. Unimodular transformations of double loops. In Proceedings of the Third Workshop on Programming Languages and Compilers for Parallel Computing, Irvine, CA. Pitman/MIT Press, 1990.
  - M.E. Wolf and M.S. Lam. A data locality optimizing algorithm. In Proceedings of the ACM SIGPLAN '91 Conference on Programming Language Design and Implementation, volume 26, pages 30-44, Toronto, Ontario, Canada, June 1991.

# Additional material for background

# A little history...early days at Bell Labs

- 1940: Russell Ohl develops PN junction (accidentally...)

- 1945: Shockley's lab established

- 1947: Bardeen and Brattain create point-contact transistor with two PN junctions, gain=18

- 1951: Shockley develops junction transistor which can be manufactured in quantity

- 1952: British radar expert GWA Dummer forecasts "solid block [with] layers of insulating, conducting and amplifying materials

- 1954: first transistor radio. Also Texas Instruments makes first silicon transistor (price $2.50)

**First point-contact transistor invented at Bell Labs.**
(Source: Bell Labs.)



**The three inventors of the transistor: William Shockley, (seated), John Bardeen (left) and Walter Brattain (right) in 1948; the three inventors shared the Nobel prize in 1956.**
(Source: Bell Labs.)

**This background section is not covered in the lectures**

# Pre-historic integrated circuits

1958: The first monolithic integrated circuit, about the size of a finger tip, developed at Texas Instruments by Jack Kilby. The IC was a chip of a single Germanium crystal containing one transistor, one capacitor, and one resistor (Source: Texas Instruments)
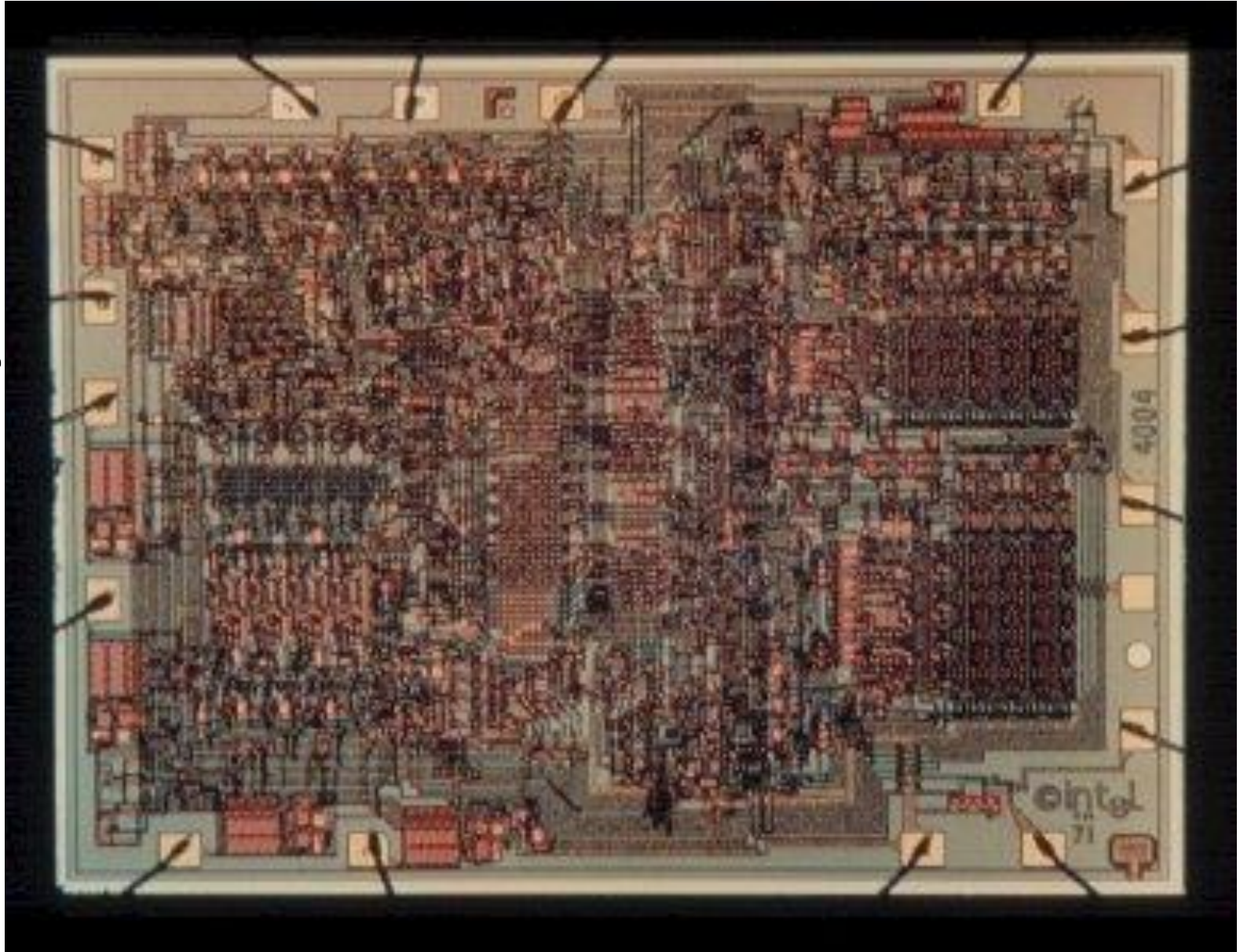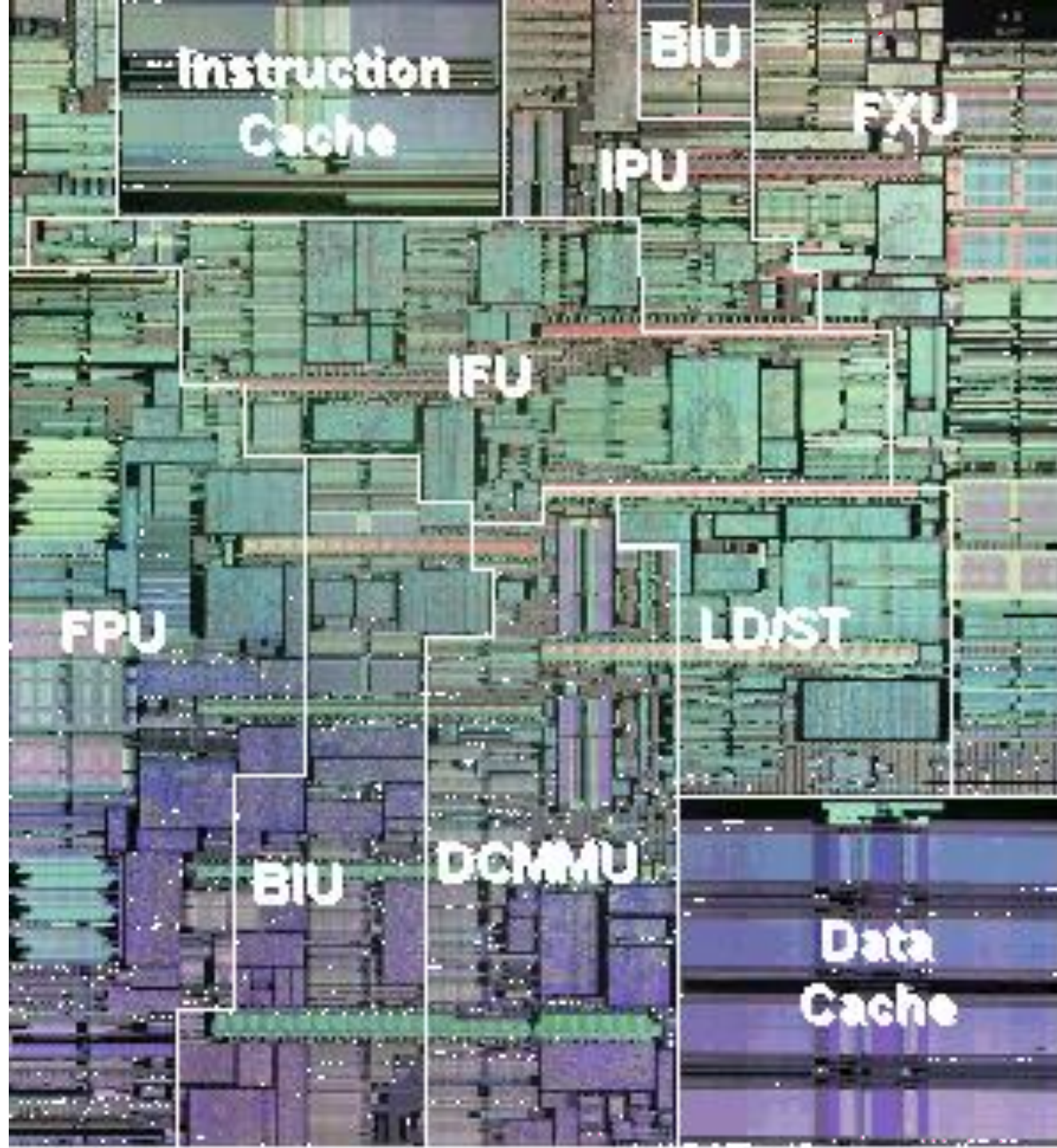
- 1970: Intel starts selling a 1K bit RAM

- 1971: Intel introduces first microprocessor, the 4004
  - 4-bit buses
  - Clock rate 108 KHz
  - 2300 transistors
  - 10$\mu$m process
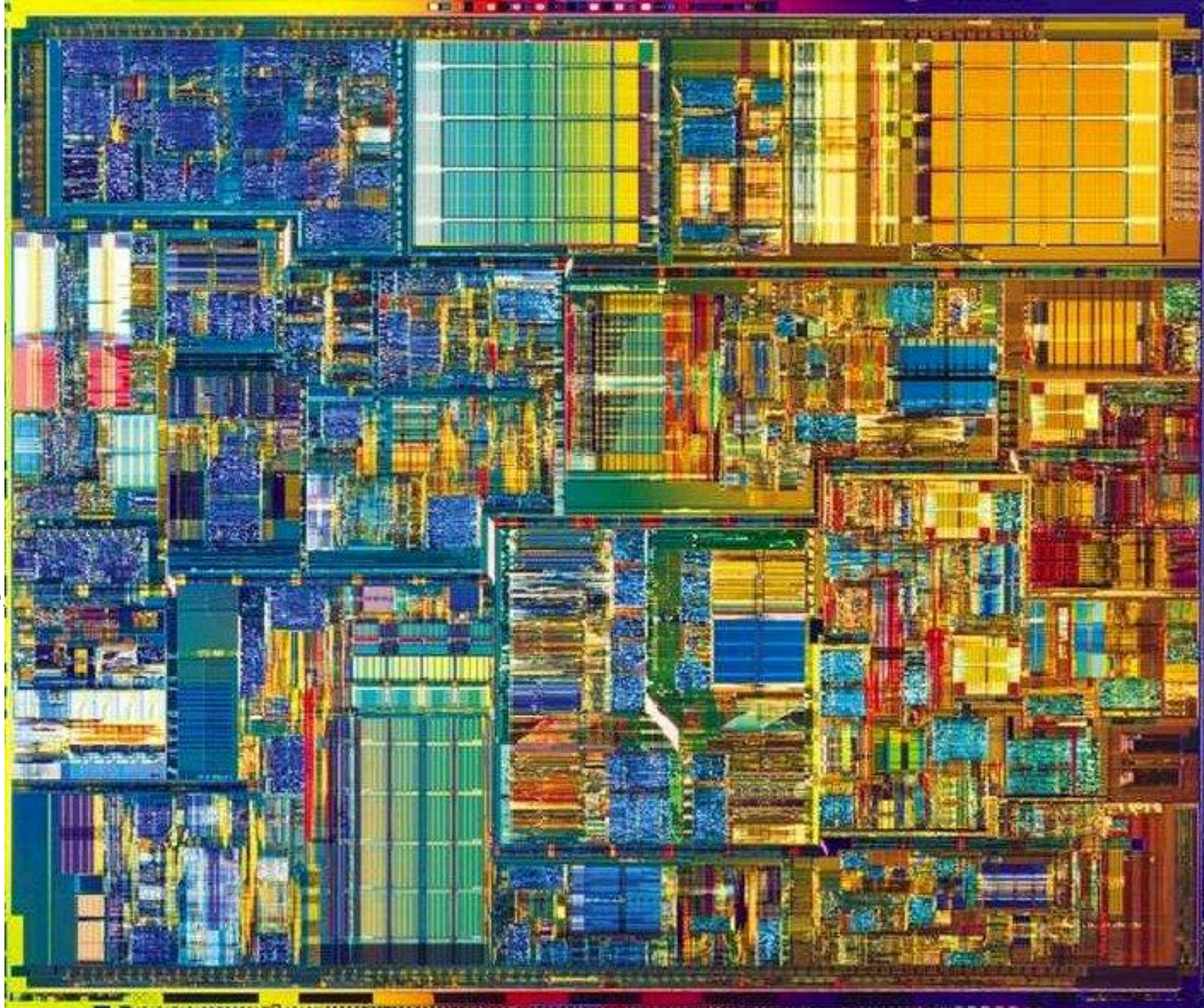
- IBM Power3 microprocessor
- 15M transistors
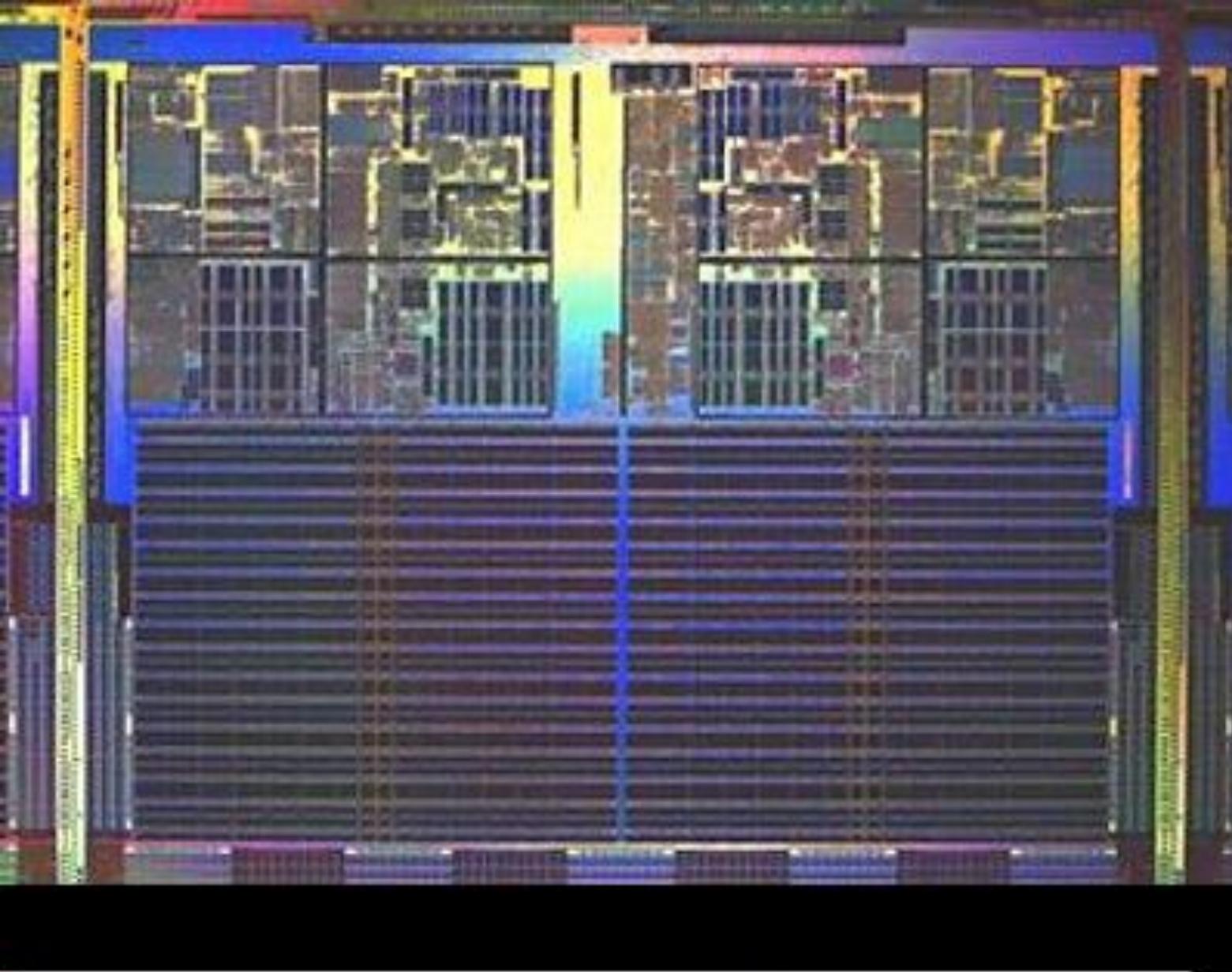- 0.18$\mu$m copper/SOI process
- About 270mm$^2$

# Intel Pentium 4

- 42 M transistors
- 0.13mm copper/SOI process
- Clock speeds: 2200, 2000MHz
- Die size 146 square mm
- Power consumption 55.1W (2200), 52.4W (2000)
- Price ($ per chip in 1,000-chip units, Jan 2002):
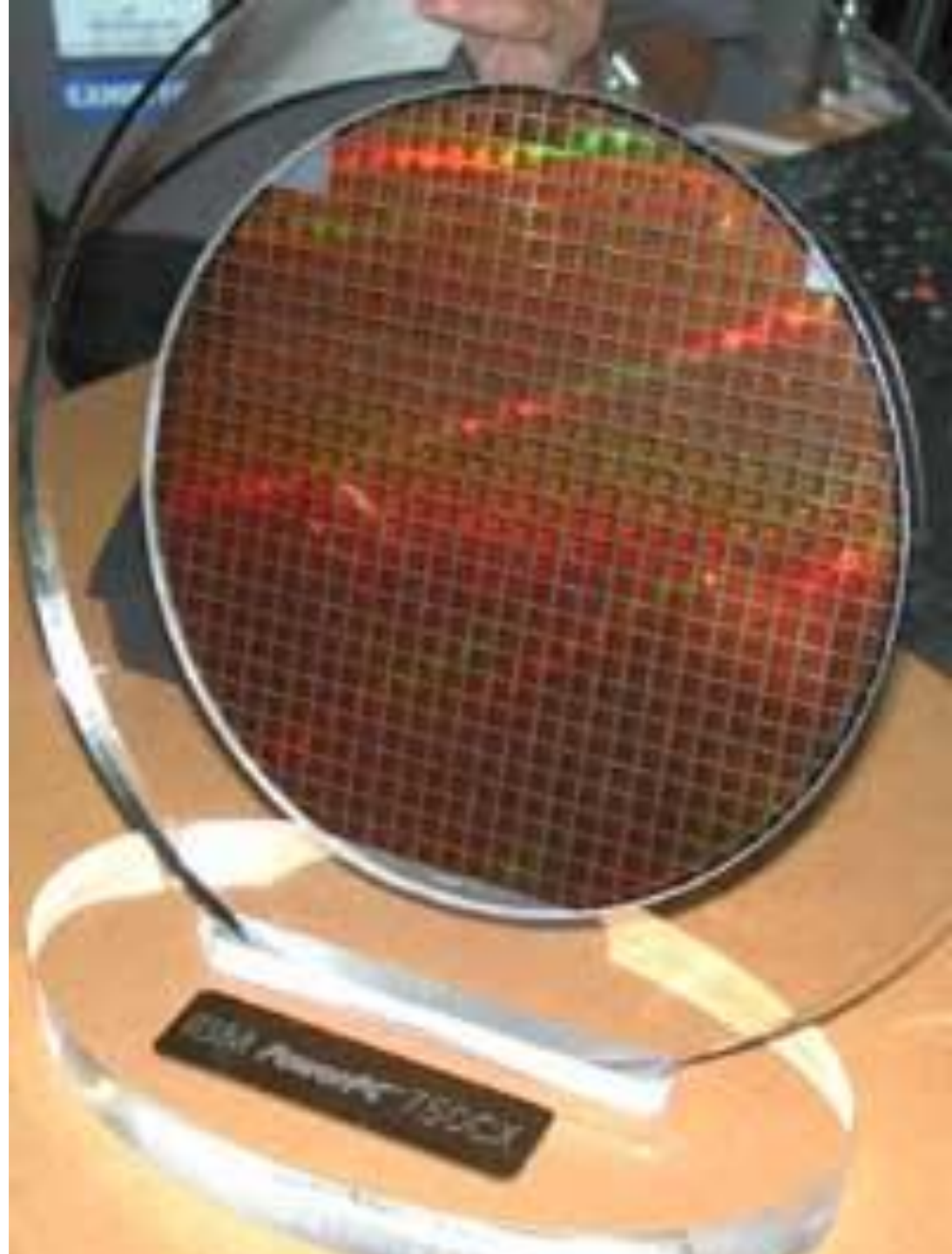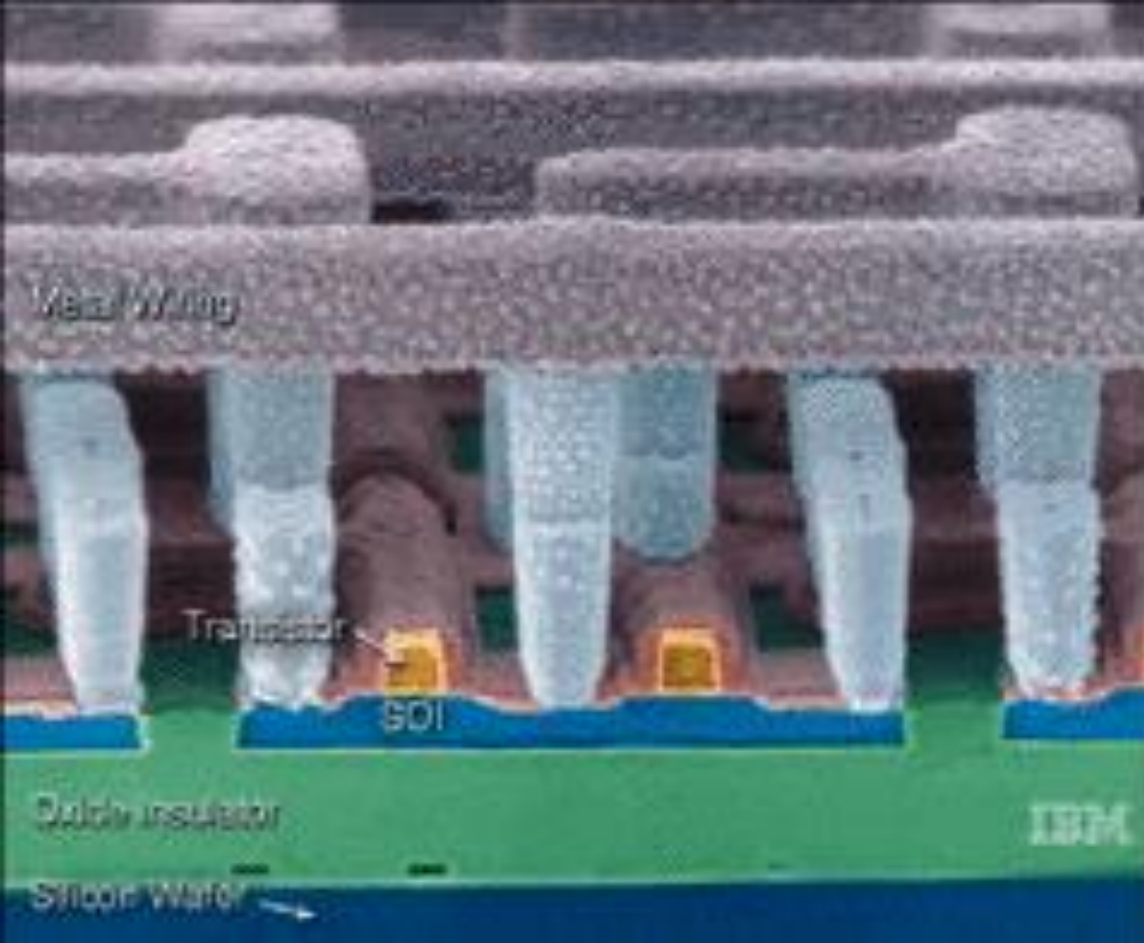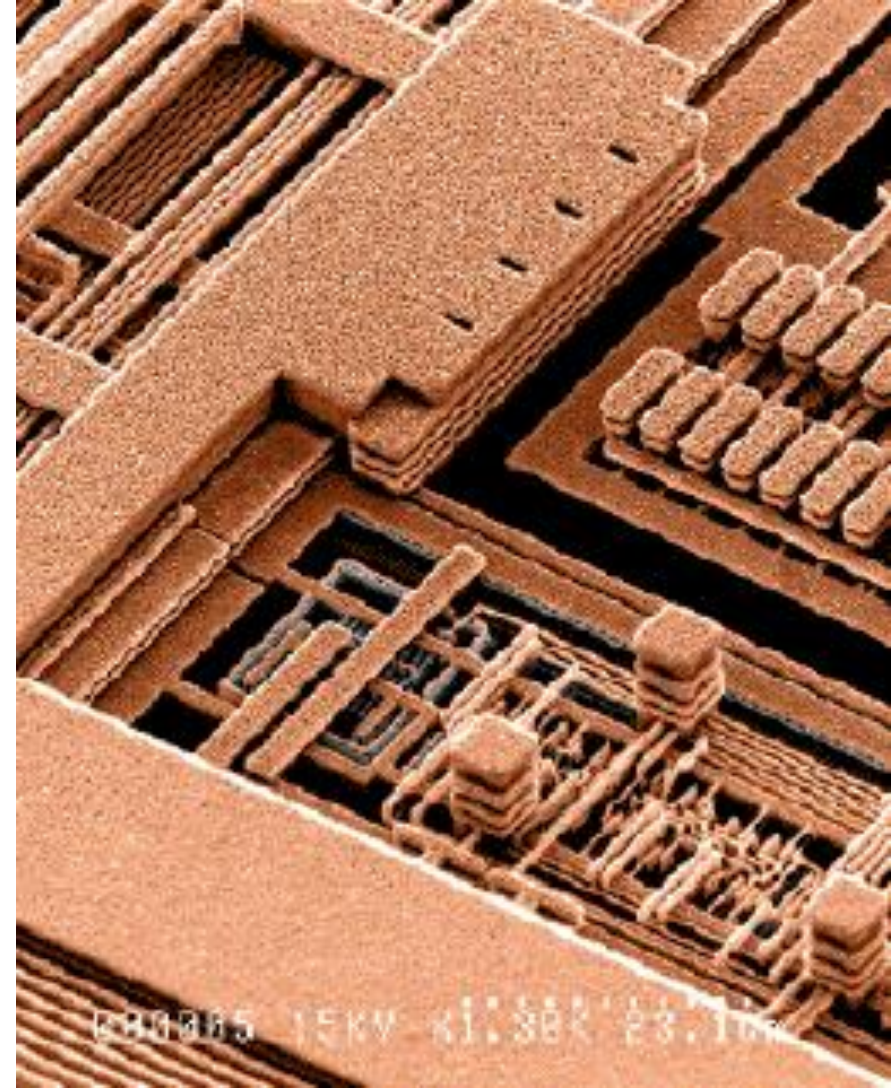  - US$562 (2200)
  - US$364 (2000)

# Dual-core Opteron 275

- **233M transistors**
- 199mm$^2$
- **90nm**
- **2.2GHz**

# Intel Itanium Montecito



Technology: 90nm bulk, 7 layers Cu
- 1.72B transistors
- 596mm²
- 2.0+GHz operation at self-selected voltage
- 100W electrical and thermal power limit
- Two 11 issue, 2 way TMT EPIC cores
- 3 level on-chip cache per core – 16K L1I, 16K L1D, 1MB L2I, 256K L2D, 12MB unified L3

21.5 mm

Labels on die:
- Branch Unit
- Floating Point
- Bus I/O
- 1MB L1I Cache
- 16KB L0I Cache
- Instr. Fetch
- Pipeline Control
- ALAT
- 16KB L0D Cache
- Integer Datapath
- HPW
- 256KB L1D Cache
- Bus Logic
- Fuse
- L2 Tag
- 12MB L2 Cache
- L2 Tag
- Bus Arbiter
- CLK
- Bus Logic
- Foxton Power management
- 12MB L2 Cache
- L1D Cache
- Integer Datapath
- HPW
- L0D
- L0I Cache ALAT
- Instr. Fetch
- L1I
- Floating Point
- Branch Unit
- Bus I/O

- Dual core
- 2-way SMT
- Launched July 2006
- List price $3692 (in 1000s)

- Chips are made from slices of a single-crystal silicon ingot
- Each slice is about 30cm in diameter, and 250-600 microns thick
- Transistors and wiring are constructed by photolithography
- Essentially a printing/etching process
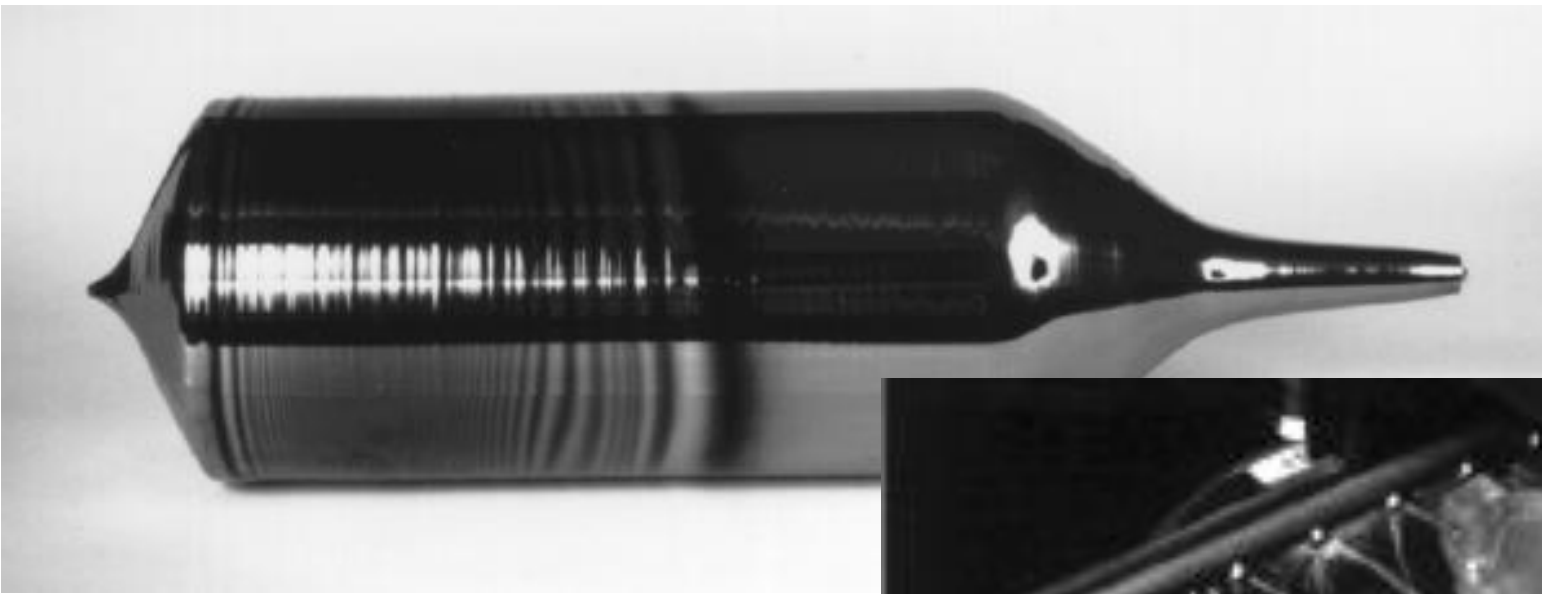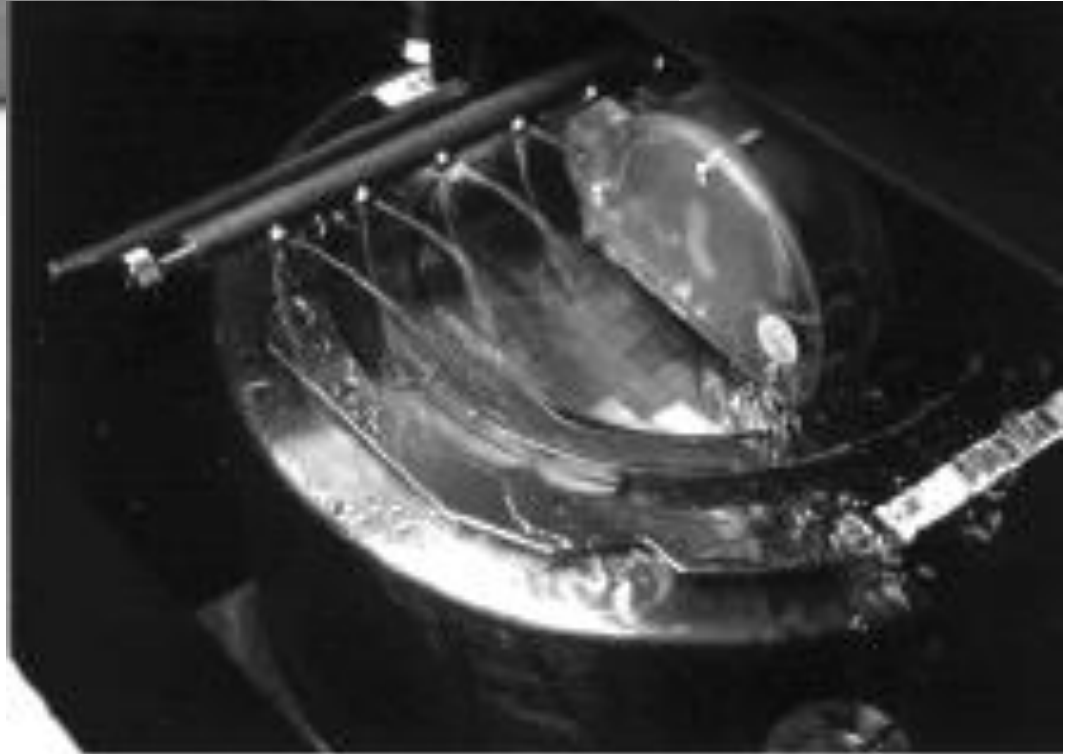- With lines ca. 0.045-0.18$\mu$m wide

> **Highly magnified scanning electron microscope (SEM) view of IBM's Silicon-On-Insulator (SOI) fabrication: a layer of silicon crystal is grown on top of a layer of insulating silicon oxide**

Highly magnified scanning electron microscope (SEM) view of IBM's six-level copper interconnect technology in an integrated circuit chip. The aluminum in transistor interconnections in a silicon chip has been replaced by copper that has a higher conductivity (by nearly 40%) and also a better ability to carry higher current densities without electromigration. Lower copper interconnect resistance means higher speeds and lower RC constants (Photograph courtesy of IBM Corporation, 1997.)
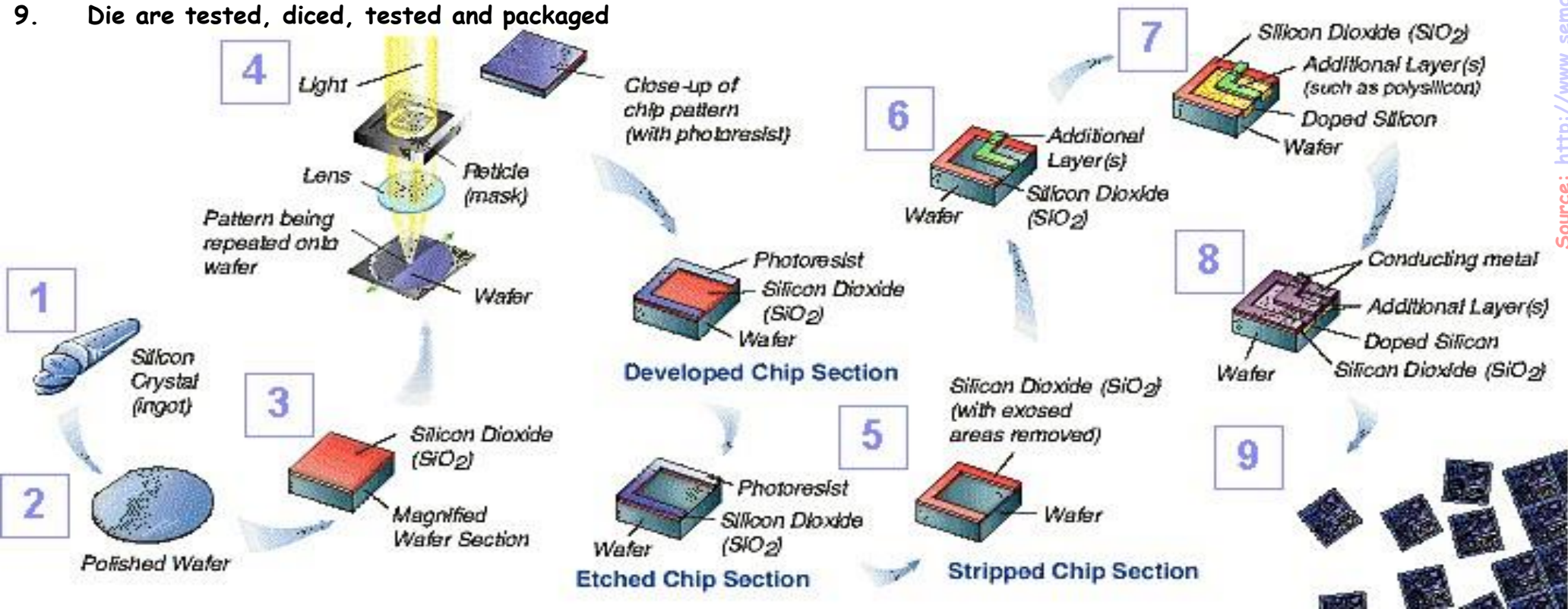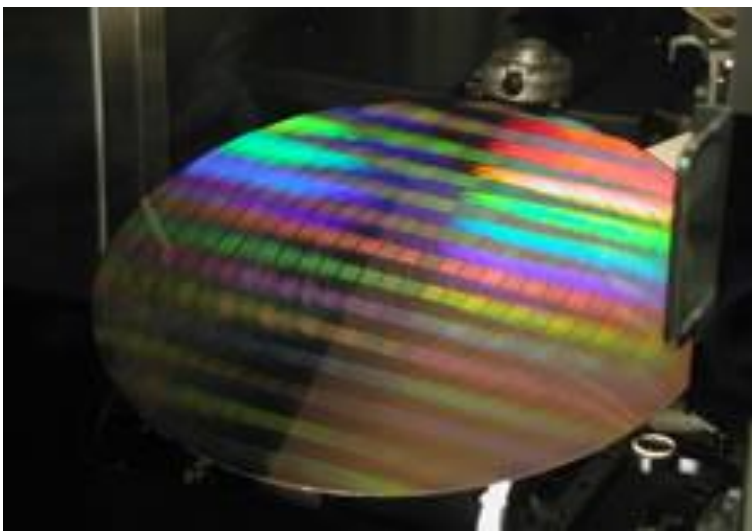
A single crystal of silicon, a silicon ingot, grown by the Czochralski technique. The diameter of this ingot is 6 inches (Courtesy of Texas Instruments).  State of the art fabs now use 300mm wafers

# Integrated circuit fabrication is a printing process

1. Grow pure silicon crystal

2. Slice into wafers and polish

3. Grow surface layer of silicon dioxide (ie glass), either using high-temperature oxygen or chemical vapour deposition

4. Coat surface with photoresist layer, then use mask to selectively expose photoresist to ultraviolet light

5. Etch away silicon dioxide regions not covered by hardened photoresist

6. Further photolithography steps build up additional layers, such as polysilicon

7. Exposed silicon is doped with small quantities of chemicals which alter its semiconductor behaviour to create transistors

8. Further photolithography steps build layers of metal for wiring

9. Die are tested, diced, tested and packaged

**4** Light — Lens — Reticle (mask) — Pattern being repeated onto wafer — Wafer

Close-up of chip pattern (with photoresist)

**6** Additional Layer(s) — Silicon Dioxide (SiO₂) — Wafer

**7** Silicon Dioxide (SiO₂) — Additional Layer(s) (such as polysilicon) — Doped Silicon — Wafer

**1** Silicon Crystal (ingot)

**3** Silicon Dioxide (SiO₂) — Magnified Wafer Section

Photoresist — Silicon Dioxide (SiO₂) — Wafer
**Developed Chip Section**

**8** Conducting metal — Additional Layer(s) — Doped Silicon — Wafer — Silicon Dioxide (SiO₂)

**2** Polished Wafer

**5** Silicon Dioxide (SiO₂) (with exosed areas removed) — Wafer
**Stripped Chip Section**

Photoresist — Silicon Dioxide (SiO₂) — Wafer
**Etched Chip Section**

**9**

Close up of the wafer as it spins during a testing procedure


Checking wafers processing in a vertical diffusion furnace


Intel technicians monitor wafers in an automated wet etch tool. The process cleans the wafers of any excess process chemicals or contamination.

# Intel x86/Pentium Family

| CPU | Year | Data Bus | Max. Mem. | Transistors | Clock MHz | Av. MIPS | Level-1 Caches |
|-----|------|----------|-----------|-------------|-----------|----------|----------------|
| **8086** | 1978 | 16 | 1MB | 29K | 5-10 | 0.8 | |
| **80286** | 1982 | 16 | 16MB | 134K | 8-12 | 2.7 | |
| **80386** | 1985 | 32 | 4GB | 275K | 16-33 | 6 | |
| **80486** | 1989 | 32 | 4GB | 1.2M | 25-100 | 20 | 8Kb |
| **Pentium** | 1993 | 64 | 4GB | 3.1M | 60-233 | 100 | 8K Instr + 8K Data |
| **Pentium Pro** | 1995 | 64 | 64GB | 5.5M +15.5M | 150-200 | 440 | 8K + 8K + Level2 |
| **Pentium II** | 1997 | 64 | 64GB | 7M | 266-450 | 466- | 16K+16K + L2 |
| **Pentium III** | 1999 | 64 | 64GB | 8.2M | 500-1000 | 1000- | 16K+16K + L2 |
| **Pentium IV** | 2001 | 64 | 64GB | 42M | 1300-2000 | | 8K + L2 |

**On-line manuals:** http://x86.ddj.com/intel.doc/386manuals.htm
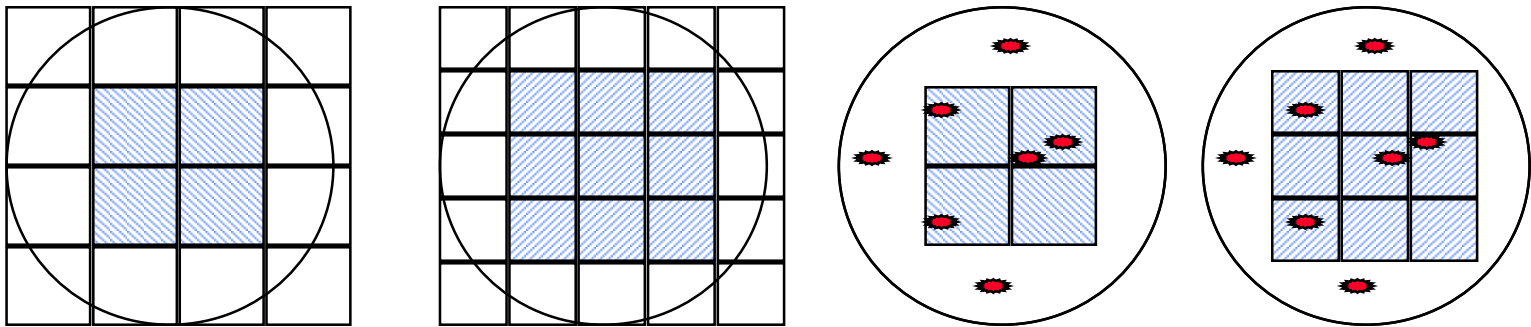
**On-line details:** http://www.sandpile.org/ia32/index.htm

# Integrated Circuits Costs

$$\text{IC cost} = \frac{\text{Die cost + Testing cost + Packaging cost}}{\text{Final test yield}}$$

$$\text{Die cost} = \frac{\text{Wafer cost}}{\text{Dies per Wafer} \times \text{Die yield}}$$

$$\text{Dies per wafer} = \frac{\pi \, (\text{Wafer\_diam}/2)^2}{\text{Die\_Area}} - \frac{\pi \times \text{Wafer\_diam}}{\sqrt{2 \cdot \text{Die\_Area}}} - \text{Test\_Die}$$
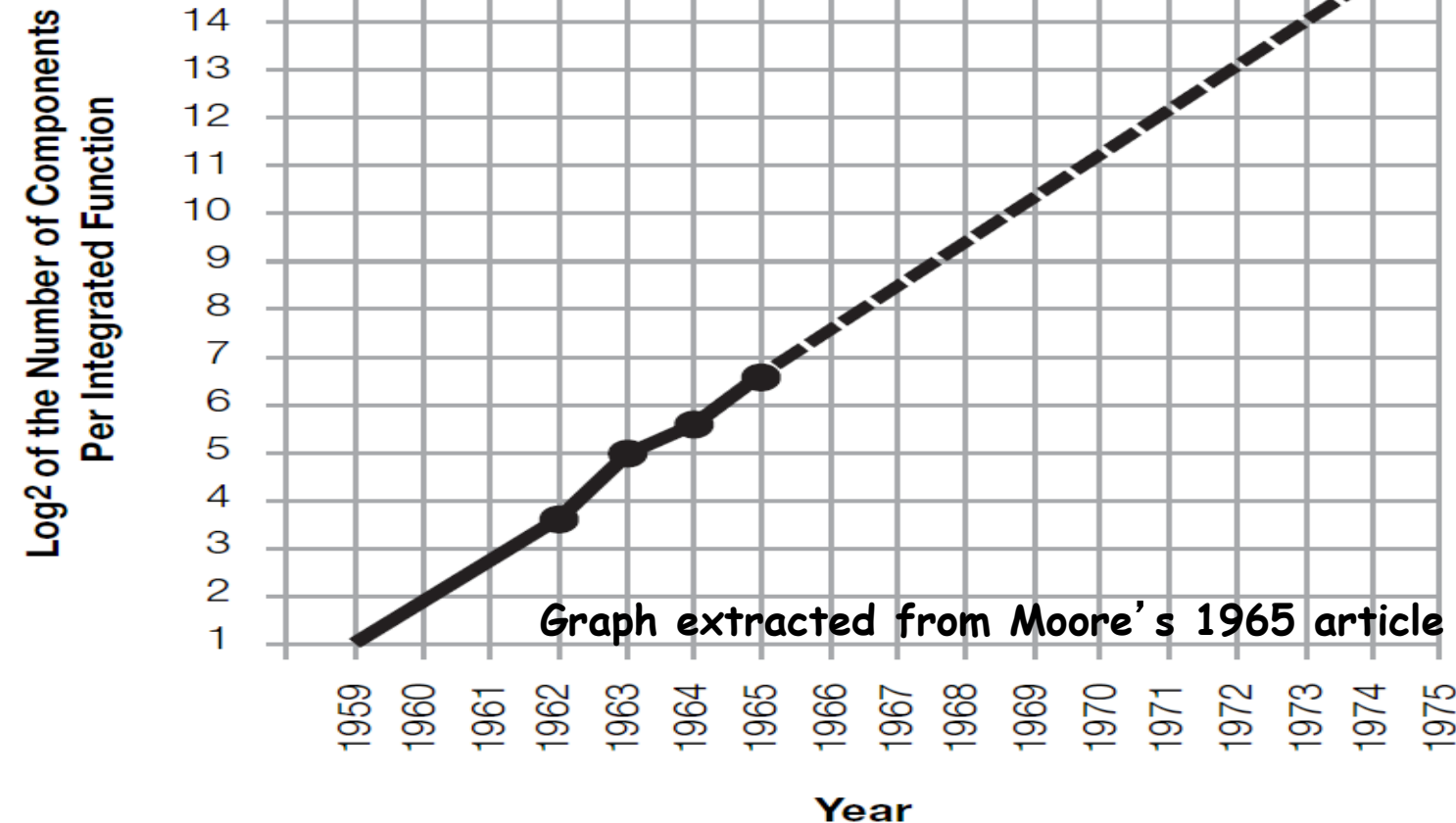
$$\text{Die Yield} = \text{Wafer\_yield} \times \left\{ 1 - \left( \frac{\text{Defect\_Density} \times \text{Die\_area}}{\alpha} \right)^{-\alpha} \right\}$$

# Die Cost goes roughly with die area$^4$

Graph extracted from Moore's 1965 article

# Moore's "Law"

Gordon Moore left Fairchild to found Intel in 1968 with Robert Noyce and Andy Grove,

**Cramming more components onto integrated circuits**
**By Gordon E. Moore**

**Electronics, Volume 38, Number 8, April 19, 1965**
**(See http://www.intel.com/research/silicon/mooreslaw.htm)**

**"With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65,000 components on a single silicon chip"**
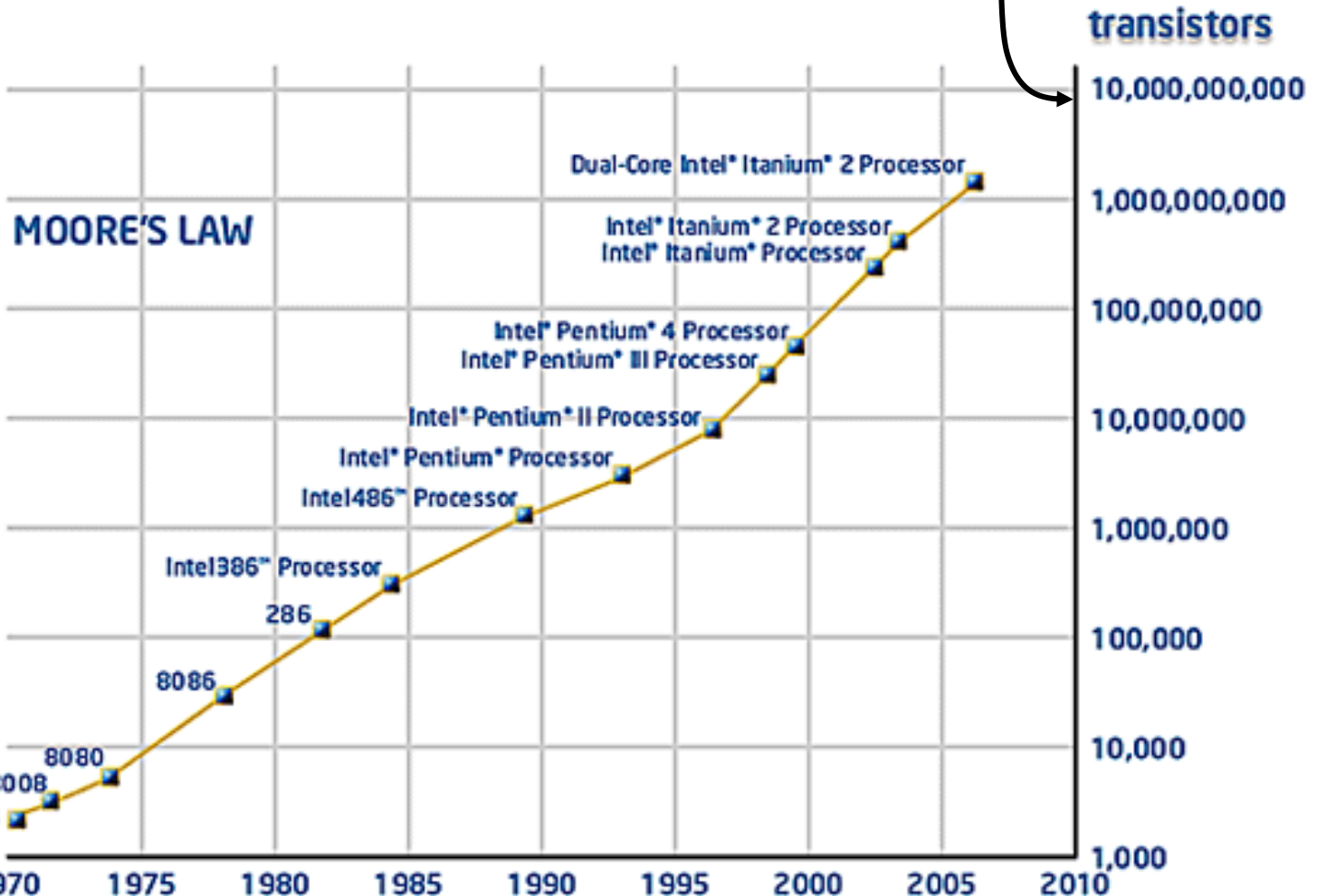
# Technology Trends: Microprocessor Capacity

CMOS improvements:
- Die size: 2X every 3 yrs
- Line width: halve / 7 yrs

"Graduation Window"

**From:**

http://www.intel.com/technology/mooreslaw/

**See also**

http://download.intel.com/research/silicon/Gordon_Moore_ISSCC_021003.pdf



transistors

MOORE'S LAW

Dual-Core Intel® Itanium® 2 Processor
Intel® Itanium® 2 Processor
Intel® Itanium® Processor
Intel® Pentium® 4 Processor
Intel® Pentium® III Processor
Intel® Pentium® II Processor
Intel® Pentium® Processor
Intel486™ Processor
Intel386™ Processor
286
8086
8080
8008
4004

# Loop-carried dependences

S1 : A[0] := 0

for I = 1 to 8

S2 : A[I] := A[I-1] + B[I]

- **Appears to be inherently sequential**

- **But parallel implementation is possible**

| B: | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
|---|---|---|---|---|---|---|---|---|---|
| >>1: | | 1 | 1 | 1 | 1 | 1 | 1 | 1 | **+** |

| B: | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | |
|---|---|---|---|---|---|---|---|---|---|
| >>2: | | | 1 | 2 | 2 | 2 | 2 | 2 | **+** |

| B: | 1 | 2 | 3 | 4 | 4 | 4 | 4 | 4 | |
|---|---|---|---|---|---|---|---|---|---|
| >>4: | | | | | 1 | 2 | 3 | 4 | **+** |

| B: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|

"Parallel scan" or "parallel prefix sum"

0   B[1]
+
   B[2]
+
   B[3]
+
   B[4]
+
   B[5]
+
   B[6]
+
   B[7]
+
   B[8]
+