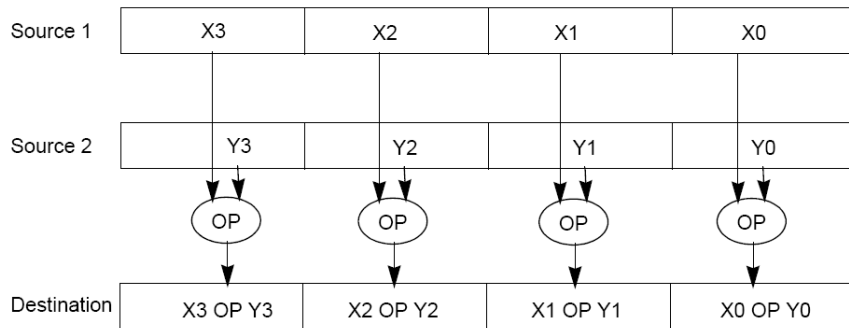# Intel SIMD extensions

1

# Performance boost

- □ Architecture improvements (such as pipeline/cache/SIMD) are more significant
- □ Intel analyzed multimedia applications and found they share the following characteristics:
  - ◻ Small native data types (8-bit pixel, 16-bit audio)
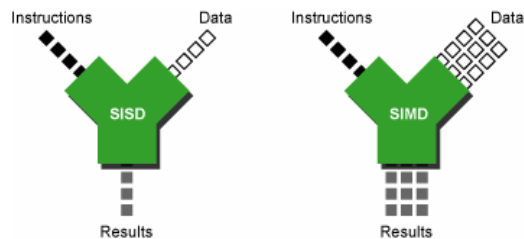  - ◻ Recurring operations
  - ◻ Inherent parallelism

2

## SIMD

- SIMD (single instruction multiple data) architecture performs the same operation on multiple data elements in parallel
- **PADDW MM0, MM1**

| Source 1 | X3 | X2 | X1 | X0 |
|---|---|---|---|---|
| Source 2 | Y3 | Y2 | Y1 | Y0 |
| | OP | OP | OP | OP |
| Destination | X3 OP Y3 | X2 OP Y2 | X1 OP Y1 | X0 OP Y0 |

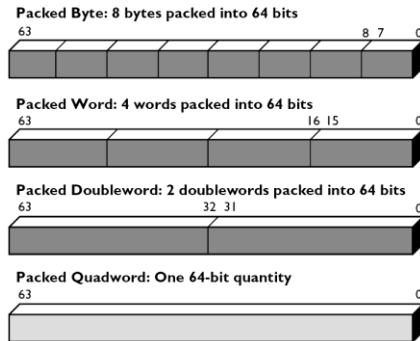## SISD/SIMD

# Intel SIMD development

- MMX (Multimedia Extension) was introduced in 1996 (Pentium with MMX and Pentium II).
- SSE (Streaming SIMD Extension) was introduced with Pentium III.
- SSE2 was introduced with Pentium 4.
- SSE3 was introduced with Pentium 4 supporting hyper-threading technology. SSE3 adds 13 more instructions.
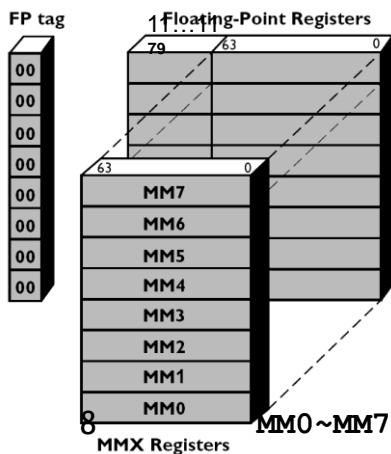- Advanced Vector Extensions (2010)

# MMX

- After analyzing a lot of existing applications such as graphics, MPEG, music, speech recognition, game, image processing, they found that many multimedia algorithms execute the same instructions on many pieces of data in a large data set.
- Typical elements are small, 8 bits for pixels, 16 bits for audio, 32 bits for graphics and general computing.
- New data type: 64-bit packed data type.

# MMX data types

Each of the MMn registers is a 64-bit integer. However, one of the main concepts of the MMX instruction set is the concept of packed data types, which means instead of using the whole register for a single 64-bit integer (quadword), two 32-bit integers (doubleword), four 16-bit integers (word) or eight 8-bit integers (byte) may be used.

Packed Byte: 8 bytes packed into 64 bits

Packed Word: 4 words packed into 64 bits

Packed Doubleword: 2 doublewords packed into 64 bits

Packed Quadword: One 64-bit quantity

# MMX integration into IA

FP tag    11...11  Floating-Point Registers
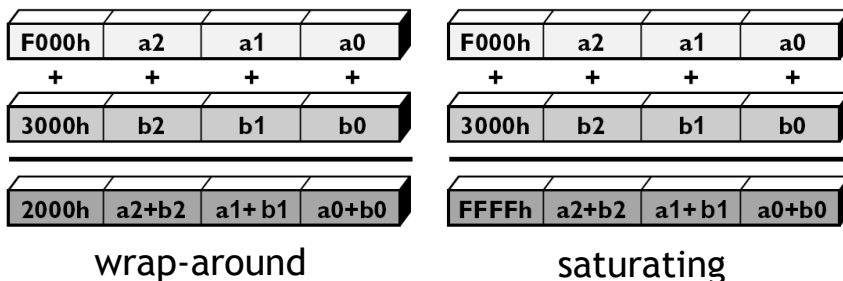
MM0~MM7

MMX Registers

- To simplify the design and to avoid modifying the operating system to preserve additional state through context switches, MMX re-uses the existing eight IA-32 FPU registers.
- This made it difficult to work with floating point and SIMD data at the same time.
- To maximize performance, programmers must use the processor exclusively in one mode or the other

# MMX instructions

- *57 MMX instructions are defined to perform the parallel operations on multiple data elements packed into 64-bit data types.*
- These include **add, subtract, multiply, compare,** and **shift, data conversion, 64-bit data move, 64-bit logical operation** and **multiply-add** for multiply-accumulate operations.
- All instructions except for data move use MMX registers as operands.
- Most complete support for 16-bit operations.

# Saturation arithmetic

- Useful in graphics applications.
- When an operation overflows or underflows, the result becomes the largest or smallest possible representable number.
- Two types: signed and unsigned saturation

| F000h | a2 | a1 | a0 |
|-------|-----|-----|-----|
| + | + | + | + |
| 3000h | b2 | b1 | b0 |
| 2000h | a2+b2 | a1+ b1 | a0+b0 |

wrap-around

| F000h | a2 | a1 | a0 |
|-------|-----|-----|-----|
| + | + | + | + |
| 3000h | b2 | b1 | b0 |
| FFFFh | a2+b2 | a1+b1 | a0+b0 |

saturating

# MMX instructions

| | Category | Wraparound | Signed Saturation | Unsigned Saturation |
|---|---|---|---|---|
| Arithmetic | Addition | PADDB, PADDW, PADDD | PADDSB, PADDSW | PADDUSB, PADDUSW |
| | Subtraction | PSUBB, PSUBW, PSUBD | PSUBSB, PSUBSW | PSUBUSB, PSUBUSW |
| | Multiplication Multiply and Add | PMULL, PMULH PMADD | | |
| Comparison | Compare for Equal | PCMPEQB, PCMPEQW, PCMPEQD | | |
| | Compare for Greater Than | PCMPGTPB, PCMPGTPW, PCMPGTPD | | |
| Conversion | Pack | | PACKSSWB, PACKSSDW | PACKUSWB |
| Unpack | Unpack High | PUNPCKHBW, PUNPCKHWD, PUNPCKHDQ | | |
| | Unpack Low | PUNPCKLBW, PUNPCKLWD, PUNPCKLDQ | | |

# MMX instructions

| | | Packed | | Full Quadword | |
|---|---|---|---|---|---|
| Logical | And And Not Or Exclusive OR | | | PAND PANDN POR PXOR | |
| Shift | Shift Left Logical Shift Right Logical Shift Right Arithmetic | PSLLW, PSLLD PSRLW, PSRLD PSRAW, PSRAD | | PSLLQ PSRLQ | |
| | | Doubleword Transfers | | Quadword Transfers | |
| Data Transfer | Register to Register Load from Memory Store to Memory | MOVD MOVD MOVD | | MOVQ MOVQ MOVQ | |
| Empty MMX State | | EMMS | | | |

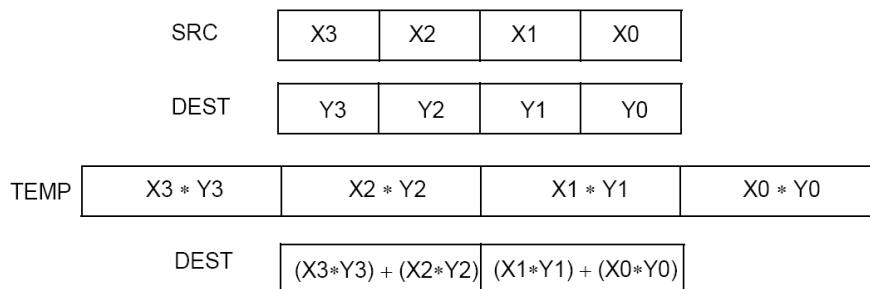Call it before you switch to FPU from MMX; Expensive operation

12

## Arithmetic

- **PADDB/PADDW/PADDD**: add two packed numbers
- Multiplication: two steps
- **PMULLW**: multiplies four words and stores the four lo words of the four double word results
- **PMULHW/PMULHUW**: multiplies four words and stores the four hi words of the four double word results. **PMULHUW** for unsigned.

13

## Arithmetic

- **PMADDWD mmi, mmj**

DEST[31:0] ← (DEST[15:0] * SRC[15:0]) + (DEST[31:16] * SRC[31:16]);
DEST[63:32] ← (DEST[47:32] * SRC[47:32]) + (DEST[63:48] * SRC[63:48]);

| | | | | |
|---|---|---|---|---|
| SRC | X3 | X2 | X1 | X0 |
| DEST | Y3 | Y2 | Y1 | Y0 |

| TEMP | X3 * Y3 | X2 * Y2 | X1 * Y1 | X0 * Y0 |
|---|---|---|---|---|

| DEST | (X3*Y3) + (X2*Y2) | (X1*Y1) + (X0*Y0) |
|---|---|---|

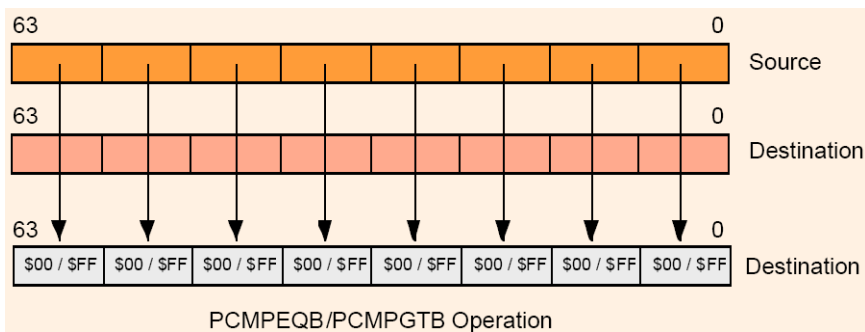## Example: add a constant to a vector

```
char d[]={5, 5, 5, 5, 5, 5, 5, 5};
char clr[]={65,66,68,...,87,88}; // 24 bytes
__asm{
    movq mm1, d
    mov cx, 3
    mov esi, 0
L1: movq mm0, clr[esi]
    paddb mm0, mm1
    movq clr[esi], mm0
    add esi, 8
    loop L1
    emms
}
```

15

## Comparison

- No CFLAGS, how many flags will you need? Results are stored in destination.
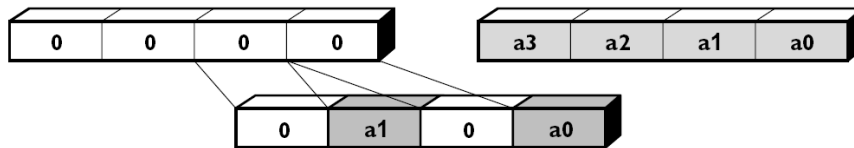- EQ/GT, no LT



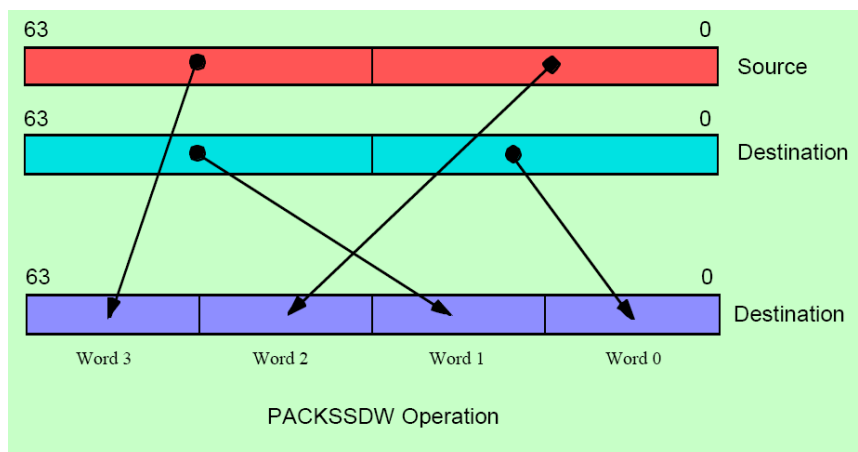PCMPEQB/PCMPGTB Operation

16

# Change data types

- □ Pack: converts a larger data type to the next smaller data type.
- □ Unpack: takes two operands and interleave them. It can be used for expand data type for immediate calculation.
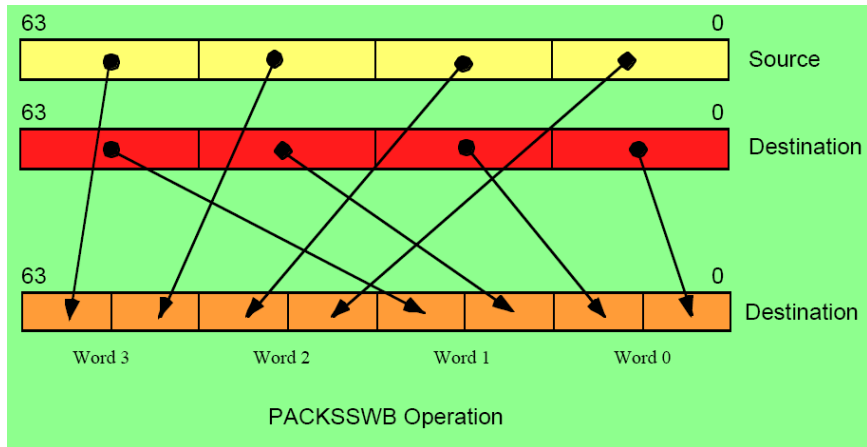
**Unpack low-order words into doublewords**

| 0 | 0 | 0 | 0 |

| a3 | a2 | a1 | a0 |

| 0 | a1 | 0 | a0 |

17

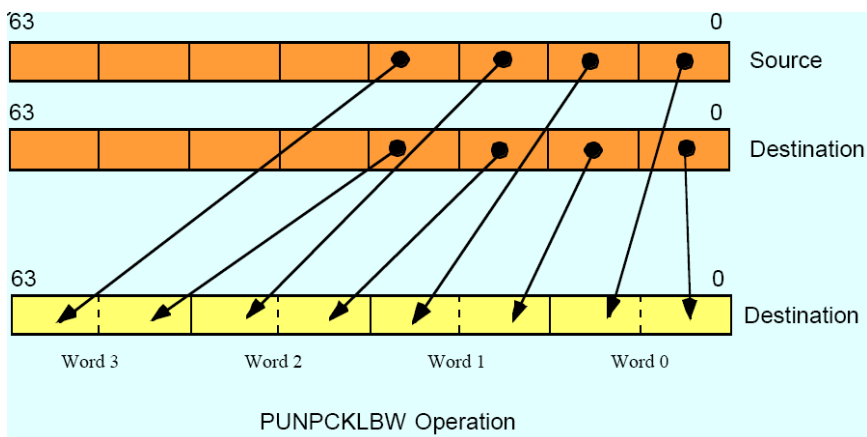# Pack with signed saturation



PACKSSDW Operation

PACKSSDW $mm_d$, $mm_s$

# Pack with signed saturation



PACKSSWB $mm_d$, $mm_s$

# Unpack low portion

# Unpack low portion



63 ... 0 — Source
63 ... 0 — Destination
63 ... 0 — Destination

DWord 1     DWord 0

PUNPCKLWD Operation

21

# Unpack low portion



63 ... 0 — Source
63 ... 0 — Destination
63 ... 0 — Destination

QWord

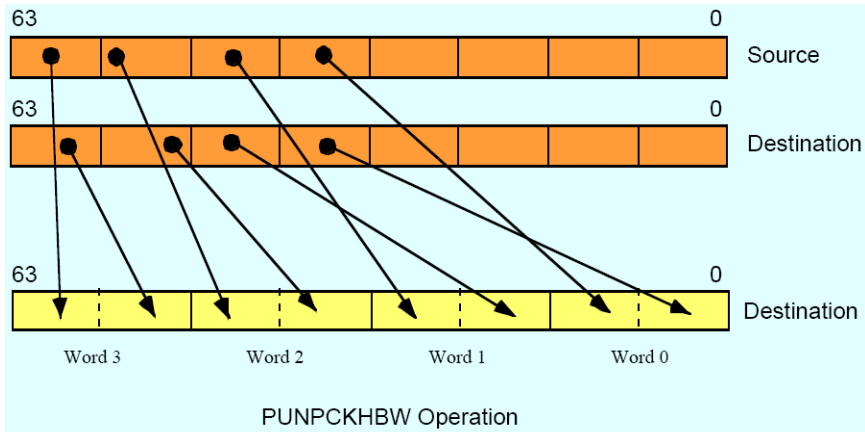PUNPCKLDQ Operation

22

11

# Unpack high portion



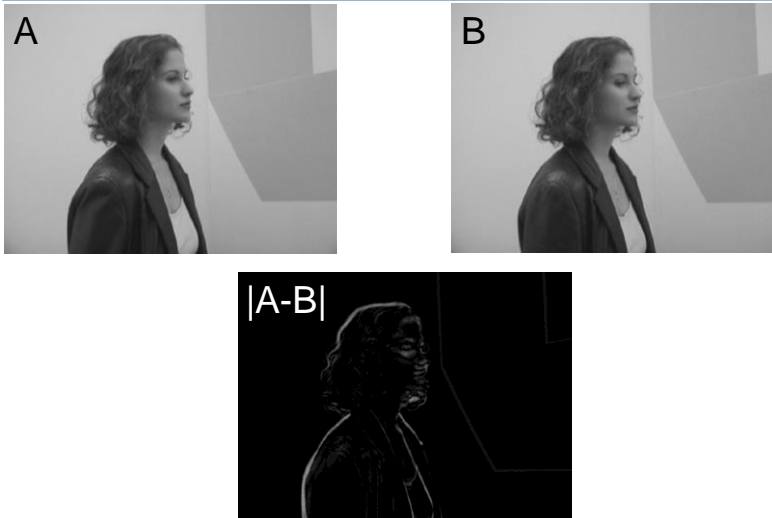PUNPCKHBW Operation

23

# Keys to SIMD programming

- Efficient data layout
- Elimination of branches

## Application: frame difference



A

B

|A-B|

## Application: frame difference



A-B

B-A

(A-B) or (B-A)

26

## Application: frame difference

```
MOVQ      mm1, A //move 8 pixels of image A
MOVQ      mm2, B //move 8 pixels of image B
MOVQ      mm3, mm1 // mm3=A
PSUBSB    mm1, mm2 // mm1=A-B
PSUBSB    mm2, mm3 // mm2=B-A
POR       mm1, mm2 // mm1=|A-B|
```

27

## Example: image fade-in-fade-out



A                          B

$A*\alpha+B*(1-\alpha) = B+\alpha(A-B)$

14

## α=0.75



## α=0.5

# α=0.25



# Example: image fade-in-fade-out

- Two formats: planar and chunky
- In Chunky format, 16 bits of 64 bits are wasted
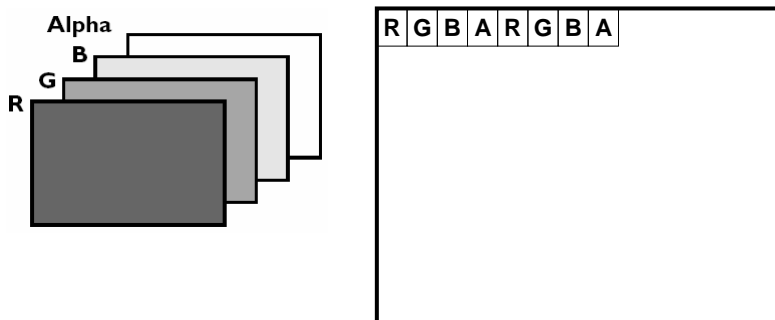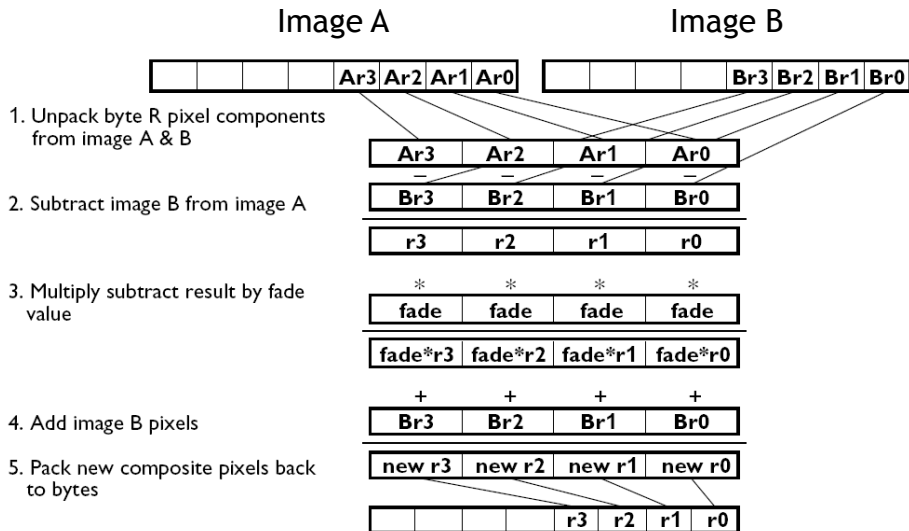- So, we use planar in the following example

# Example: image fade-in-fade-out

Image A                          Image B

| | | | | Ar3 | Ar2 | Ar1 | Ar0 | | | | | | Br3 | Br2 | Br1 | Br0 |

1. Unpack byte R pixel components
   from image A & B

| Ar3 | Ar2 | Ar1 | Ar0 |
| --- | --- | --- | --- |
| − | − | − | − |

2. Subtract image B from image A

| Br3 | Br2 | Br1 | Br0 |
| --- | --- | --- | --- |

| r3 | r2 | r1 | r0 |
| --- | --- | --- | --- |

3. Multiply subtract result by fade value

| * | * | * | * |
| --- | --- | --- | --- |
| fade | fade | fade | fade |

| fade*r3 | fade*r2 | fade*r1 | fade*r0 |
| --- | --- | --- | --- |

| + | + | + | + |
| --- | --- | --- | --- |

4. Add image B pixels

| Br3 | Br2 | Br1 | Br0 |
| --- | --- | --- | --- |

5. Pack new composite pixels back to bytes

| new r3 | new r2 | new r1 | new r0 |
| --- | --- | --- | --- |

| | | | | | r3 | r2 | r1 | r0 |

# Example: image fade-in-fade-out

```
MOVQ       mm0, alpha//4 16-b zero-padding α
MOVD       mm1, A //move 4 pixels of image A
MOVD       mm2, B //move 4 pixels of image B
PXOR       mm3, mm3 //clear mm3 to all zeroes
//unpack 4 pixels to 4 words
PUNPCKLBW mm1, mm3 // Because B-A could be
PUNPCKLBW mm2, mm3 // negative, need 16 bits
PSUBW      mm1, mm2 //(B-A)
PMULHW     mm1, mm0 //(B-A)*fade/256
PADDW      mm1, mm2 //(B-A)*fade + B
//pack four words back to four bytes
PACKUSWB   mm1, mm3
```
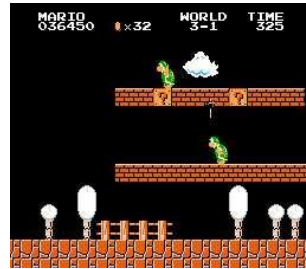
# Data-independent computation

- Each operation can execute without needing to know the results of a previous operation.
- Example, sprite overlay

```
for i=1 to sprite_Size
    if  sprite[i]=clr
    then out_color[i]=bg[i]
    else out_color[i]=sprite[i]
```



- How to execute data-dependent calculations on several pixels in parallel.

35

# Application: sprite overlay

| **Phase 1** | a3 | a2 | a1 | a0 |
|---|---|---|---|---|
| | = | = | = | = |
| | clear_color | clear_color | clear_color | clear_color |

| 1111…1111 | 0000…0000 | 1111…1111 | 0000…0000 |
|---|---|---|---|

**Phase 2**

| a3 | a2 | a1 | a0 | | c3 | c2 | c1 | c0 |
|---|---|---|---|---|---|---|---|---|
| **A and (Complement of Mask)** | | | | | **C and Mask** | | | |
| 0000…0000 | 1111…1111 | 0000…0000 | 1111…1111 | | 1111…1111 | 0000…0000 | 1111…1111 | 0000…0000 |

| 0 | a2 | 0 | a0 | | c3 | 0 | c1 | 0 |
|---|---|---|---|---|---|---|---|---|

**OR the two results
to finish the overlay**

| c3 | a2 | c1 | a0 |
|---|---|---|---|

## Application: sprite overlay

```
MOVQ     mm0, sprite
MOVQ     mm2, mm0
MOVQ     mm4, bg
MOVQ     mm1, clr
PCMPEQW  mm0, mm1
PAND     mm4, mm0
PANDN    mm0, mm2
POR      mm0, mm4
```
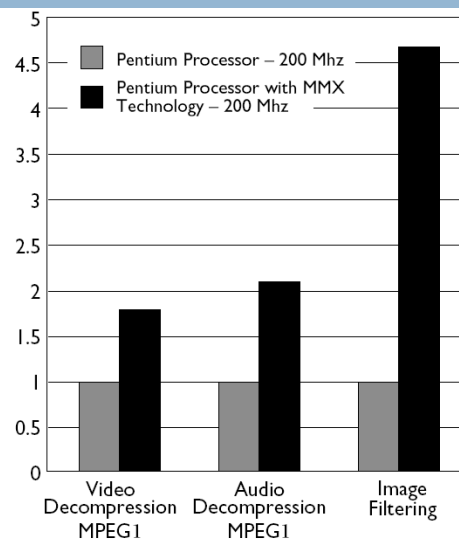
37

## Performance boost (data from 1996)

Benchmark kernels: FFT, FIR, vector dot-product, IDCT, motion compensation.

65% performance gain

Lower the cost of multimedia programs by removing the need of specialized DSP chips
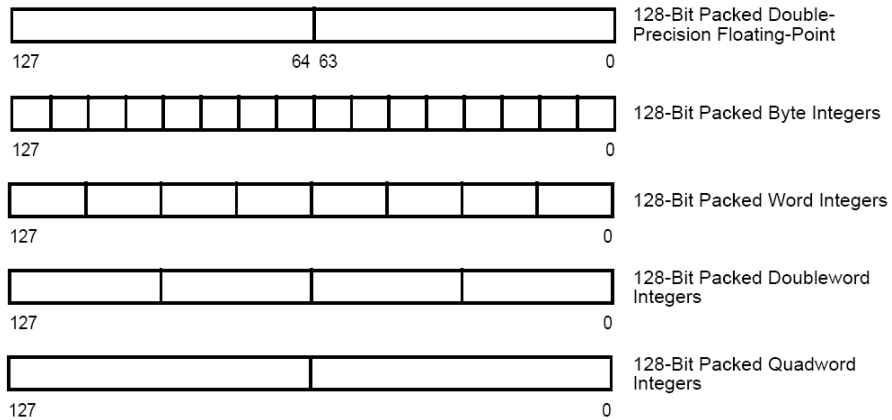
## SSE

- Adds eight 128-bit registers
- Allows SIMD operations on packed single-precision floating-point numbers
- Most SSE instructions require 16-aligned addresses

39

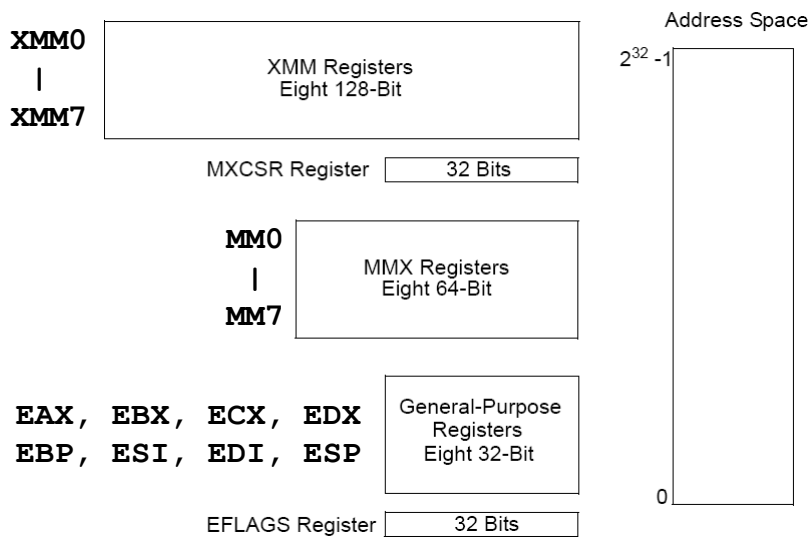## SSE features

- Add eight 128-bit data registers (XMM registers) in non-64-bit modes; sixteen XMM registers are available in 64-bit mode.
- 32-bit MXCSR register (control and status)
- Add a new data type: 128-bit packed single-precision floating-point (4 FP numbers.)
- Instruction to perform SIMD operations on 128-bit packed single-precision FP and additional 64-bit SIMD integer operations.

## SSE2 features

| | |
|---|---|
| 127      64 63      0 | 128-Bit Packed Double-Precision Floating-Point |
| 127      0 | 128-Bit Packed Byte Integers |
| 127      0 | 128-Bit Packed Word Integers |
| 127      0 | 128-Bit Packed Doubleword Integers |
| 127      0 | 128-Bit Packed Quadword Integers |

41

## SSE programming environment

**XMM0** | **XMM7**
XMM Registers Eight 128-Bit

MXCSR Register — 32 Bits

**MM0** | **MM7**
MMX Registers Eight 64-Bit

**EAX, EBX, ECX, EDX**
**EBP, ESI, EDI, ESP**
General-Purpose Registers Eight 32-Bit

EFLAGS Register — 32 Bits

Address Space
$2^{32}-1$

0

# SSE packed FP operation

| X3 | X2 | X1 | X0 |
|---|---|---|---|

| Y3 | Y2 | Y1 | Y0 |
|---|---|---|---|

OP    OP    OP    OP

| X3 OP Y3 | X2 OP Y2 | X1 OP Y1 | X0 OP Y0 |
|---|---|---|---|

☐ **ADDPS/SUBPS**: packed single-precision FP

43

# SSE scalar FP operation

| X3 | X2 | X1 | X0 |
|---|---|---|---|

| Y3 | Y2 | Y1 | Y0 |
|---|---|---|---|

OP

| X3 | X2 | X1 | X0 OP Y0 |
|---|---|---|---|

- **ADDSS/SUBSS**: scalar single-precision FP
  used as FPU?

44

## SSE2

- Provides ability to perform SIMD operations on double-precision FP, allowing advanced graphics such as ray tracing
- Provides greater throughput by operating on 128-bit packed integers

45

## Example

```
void add(float *a, float *b, float *c) {
  for (int i = 0; i < 4; i++)
    c[i] = a[i] + b[i];
}
__asm {
mov    eax, a
mov    edx, b
mov    ecx, c
movaps xmm0, XMMWORD PTR [eax]
addps  xmm0, XMMWORD PTR [edx]
movaps XMMWORD PTR [ecx], xmm0
}
```
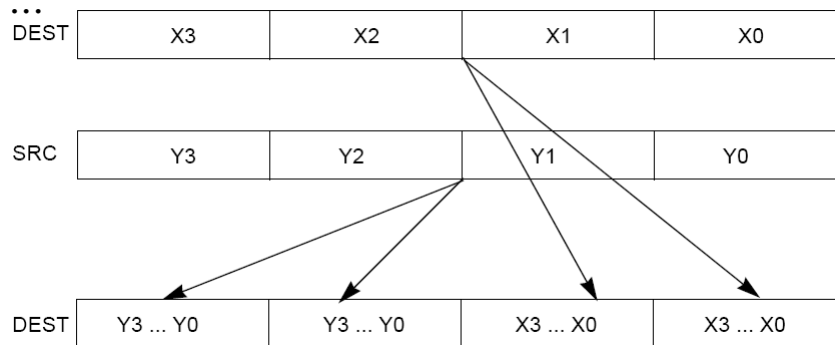
movaps: move aligned packed single-precision FP
addps: add packed single-precision FP

46

## SSE Shuffle (`SHUFPS`)

**`SHUFPS xmm1, xmm2, imm8`**

Select[1..0] decides which DW of DEST to be copied to the
  1st DW of DEST

...

| DEST | X3 | X2 | X1 | X0 |
|------|----|----|----|----|

| SRC | Y3 | Y2 | Y1 | Y0 |
|-----|----|----|----|----|

| DEST | Y3 ... Y0 | Y3 ... Y0 | X3 ... X0 | X3 ... X0 |
|------|-----------|-----------|-----------|-----------|

## SSE Shuffle (`SHUFPS`)

CASE (SELECT[1:0]) OF
    0:    DEST[31:0] ← DEST[31:0];
    1:    DEST[31:0] ← DEST[63:32];
    2:    DEST[31:0] ← DEST[95:64];
    3:    DEST[31:0] ← DEST[127:96];
ESAC;

CASE (SELECT[3:2]) OF
    0:    DEST[63:32] ← DEST[31:0];
    1:    DEST[63:32] ← DEST[63:32];
    2:    DEST[63:32] ← DEST[95:64];
    3:    DEST[63:32] ← DEST[127:96];
ESAC;

CASE (SELECT[5:4]) OF
    0:    DEST[95:64] ← SRC[31:0];
    1:    DEST[95:64] ← SRC[63:32];
    2:    DEST[95:64] ← SRC[95:64];
    3:    DEST[95:64] ← SRC[127:96];
ESAC;

CASE (SELECT[7:6]) OF
    0:    DEST[127:96] ← SRC[31:0];
    1:    DEST[127:96] ← SRC[63:32];
    2:    DEST[127:96] ← SRC[95:64];
    3:    DEST[127:96] ← SRC[127:96];
ESAC;

## Example (cross product)

```
Vector cross(const Vector& a , const Vector& b ) {
    return Vector(
        ( a[1] * b[2] - a[2] * b[1] ) ,
        ( a[2] * b[0] - a[0] * b[2] ) ,
        ( a[0] * b[1] - a[1] * b[0] ) );
}
```

49

## Example (cross product)

```
/* cross */
__m128 _mm_cross_ps( __m128 a , __m128 b ) {
  __m128 ea , eb;
  // set to a[1][2][0][3] , b[2][0][1][3]
  ea = _mm_shuffle_ps( a, a, _MM_SHUFFLE(3,0,2,1) );
  eb = _mm_shuffle_ps( b, b, _MM_SHUFFLE(3,1,0,2) );
  // multiply
  __m128 xa = _mm_mul_ps( ea , eb );
  // set to a[2][0][1][3] , b[1][2][0][3]
  a = _mm_shuffle_ps( a, a, _MM_SHUFFLE(3,1,0,2) );
  b = _mm_shuffle_ps( b, b, _MM_SHUFFLE(3,0,2,1) );
  // multiply
  __m128 xb = _mm_mul_ps( a , b );
  // subtract
  return _mm_sub_ps( xa , xb );
}
```
50