

# **RISC-V Vector Extension**

aka RVV

# Vectors... as we know them in math

1 vector with 4 elements

E1	E2	E3	E4
----	----	----	----

2 vectors with 4 elements

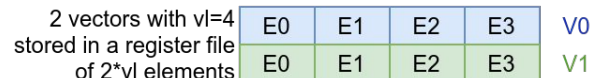
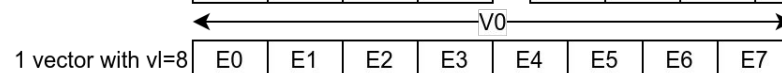
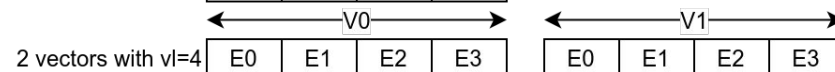
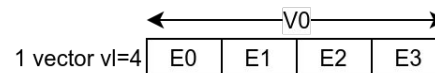
E1	E2	E3	E4	E1	E2	E3	E4
----	----	----	----	----	----	----	----

1 vector with 8 elements

E1	E2	E3	E4	E5	E6	E7	E8
----	----	----	----	----	----	----	----

2 vectors with 4 elements  
stored in a matrix 4x2

E1	E2	E3	E4
E1	E2	E3	E4

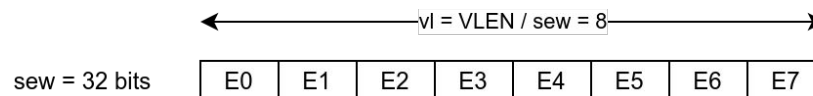
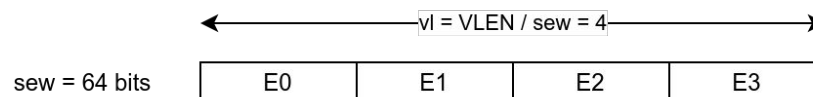
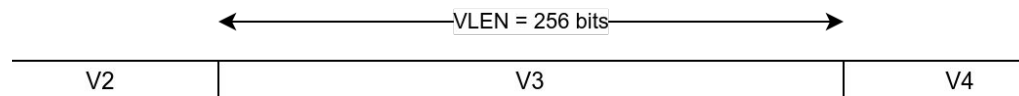


! For mapping this into hardware,  
we need to fix some parameters

# Implementation parameters

- **VLEN**, size of the vectors in bits

- Micro-architectural parameter
- RISC-V,  $VLEN \in [128; 2^{16}]$  bits
- RISC-V EPAC:  $2^{15} = 16384$  bits
- x86\_64 SSE: 128 bits
- x86\_64 AVX: 256 bits
- x86\_64 AVX512: 512 bits
- Arm SVE:  $VLEN \in [128; 2^{11}]$  bits
- Arm Neon: 128 bits



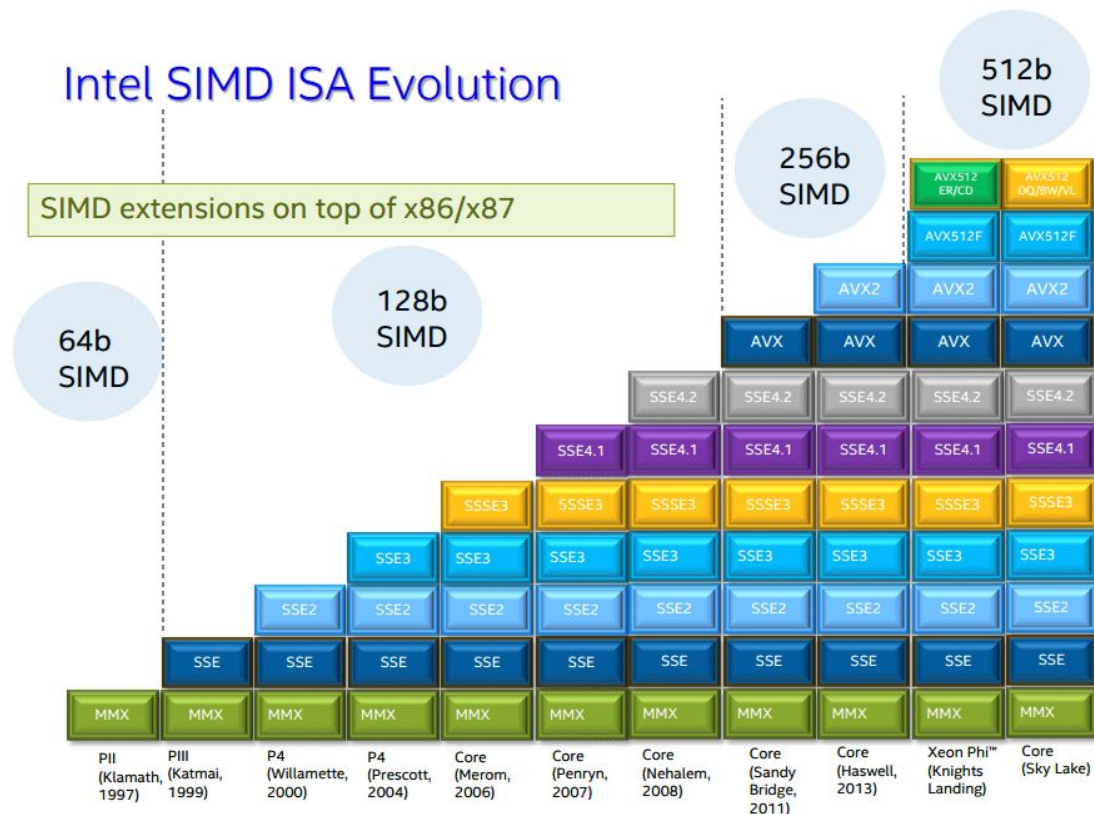
- **sew**, Selected Element Width, size of the element of a vector in bits

- You can see it as  $8 * \text{sizeof}(\text{type})$ , 64 bits for doubles/long, 32 bits for floats/int, 16 bits for (half)/short

- **VI** =  $VLEN/sew$ , the number of elements operated by the vector instruction


# An example: x86

- VLEN = 128
  - SSE
- VLEN = 256
  - AVX
- VLEN = 512
  - AVX512



# An example: AXPY with x86 intrinsics


```
43 int main()
44 {
45     int n = 32;
46     float alpha = 2.0;
47
48     // Allocate memory for vectors x and y
49     float* x = (float*)malloc(n * sizeof(float));
50     float* y = (float*)malloc(n * sizeof(float));
51
52     // Initialize vectors x and y with sample values
53     for (int i = 0; i < n; i++) {
54         x[i] = i + 1;
55         y[i] = i + 6;
56     }
57
58     printf("Original y:\n");
59     for (int i = 0; i < n; i++) {
60         printf("%f ", y[i]);
61     }
62     printf("\n");
63
64     axpy(n, alpha, x, y);
65     //axpy_avx512(n, alpha, x, y);
66     //axpy_avx512_tail(n, alpha, x, y);
67
68     printf("Resulting y after AXPY operation:\n");
69     for (int i = 0; i < n; i++) {
70         printf("%f ", y[i]);
71     }
72     printf("\n");
73
74     // Free the allocated memory
75     free(x);
76     free(y);
77
78     return 0;
79 }
```



```
36 void axpy(int n, float alpha, float *x, float *y)
37 {
38     for (int i = 0; i < n; i++) {
39         y[i] = alpha * x[i] + y[i];
40     }
41 }
```

# An example: AXPY with x86 intrinsics

```
43 int main()
44 {
45     int n = 32;
46     float alpha = 2.0;
47
48     // Allocate memory for vectors x and y
49     float* x = (float*)malloc(n * sizeof(float));
50     float* y = (float*)malloc(n * sizeof(float));
51
52     // Initialize vectors x and y with sample values
53     for (int i = 0; i < n; i++) {
54         x[i] = i + 1;
55         y[i] = i + 6;
56     }
57
58     printf("Original y:\n");
59     for (int i = 0; i < n; i++) {
60         printf("%f ", y[i]);
61     }
62     printf("\n");
63
64     axpy(n, alpha, x, y);
65     //axpy_avx512(n, alpha, x, y);
66     //axpy_avx512_tail(n, alpha, x, y);
67
68     printf("Resulting y after AXPY operation:\n");
69     for (int i = 0; i < n; i++) {
70         printf("%f ", y[i]);
71     }
72     printf("\n");
73
74     // Free the allocated memory
75     free(x);
76     free(y);
77
78     return 0;
79 }
```



```
24 void axpy_avx512(int n, float alpha, float *x, float *y)
25 {
26     int i;
27     __m512 alpha_vec = _mm512_set1_ps(alpha);
28     for (i = 0; i < n; i += 16) {
29         __m512 x_vec = _mm512_loadu_ps(&x[i]);
30         __m512 y_vec = _mm512_loadu_ps(&y[i]);
31         __m512 result = _mm512_fmadd_ps(alpha_vec, x_vec, y_vec);
32         _mm512_storeu_ps(&y[i], result);
33     }
34 }
```

# An example: AXPY with x86 intrinsics

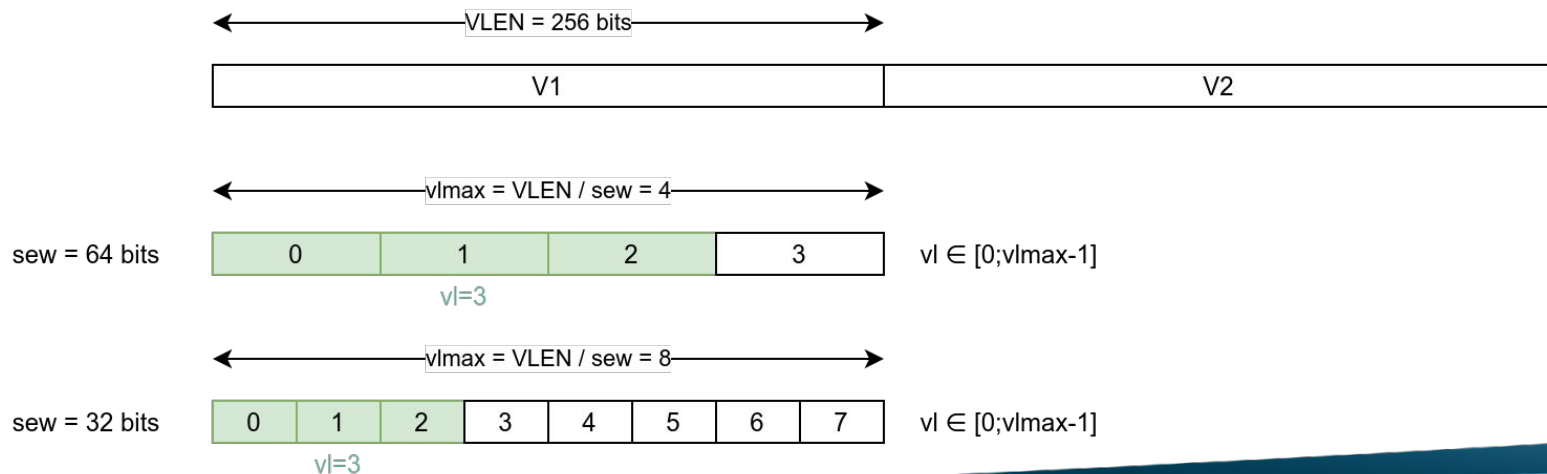
```
43 int main()
44 {
45     int n = 32;
46     float alpha = 2.0;
47
48     // Allocate memory for vectors x and y
49     float* x = (float*)malloc(n * sizeof(float));
50     float* y = (float*)malloc(n * sizeof(float));
51
52     // Initialize vectors x and y with sample values
53     for (int i = 0; i < n; i++) {
54         x[i] = i + 1;
55         y[i] = i + 6;
56     }
57
58     printf("Original y:\n");
59     for (int i = 0; i < n; i++) {
60         printf("%f ", y[i]);
61     }
62     printf("\n");
63
64     axpy(n, alpha, x, y);
65     //axpy_avx512(n, alpha, x, y);
66     //axpy_avx512_tail(n, alpha, x, y);
67
68     printf("Resulting y after AXPY operation:\n");
69     for (int i = 0; i < n; i++) {
70         printf("%f ", y[i]);
71     }
72     printf("\n");
73
74     // Free the allocated memory
75     free(x);
76     free(y);
77
78     return 0;
79 }
```

For a generic size of X and Y,  
we must handle “loop tails”

```
5 void axpy_avx512_tail(int n, float alpha, float *x, float *y)
6 {
7     int i;
8     _m512 alpha_vec = _mm512_set1_ps(alpha);
9     int avx512_loop_size = n - (n % 16);
10
11     for (i = 0; i < avx512_loop_size; i += 16) {
12         _m512 x_vec = _mm512_loadu_ps(&x[i]);
13         _m512 y_vec = _mm512_loadu_ps(&y[i]);
14         _m512 result = _mm512_fmadd_ps(alpha_vec, x_vec, y_vec);
15         _mm512_storeu_ps(&y[i], result);
16     }
17
18     for (; i < n; i++) {
19         y[i] = alpha * x[i] + y[i];
20     }
21 }
```

# A bit more flexible: variable vl

- **VLEN**, size of the vectors in bits
- **sew**, Selected Element Width, size of the element of a vector in bits
- **vlmax** =  $VLEN/sew$ , the number of elements operated by the vector instruction
- **vl**  $\in [0; vlmax-1]$ 
  - $vl = 0 \rightarrow$  yes, it's a nop!
  - $vl = vlmax \rightarrow$  back to previous case

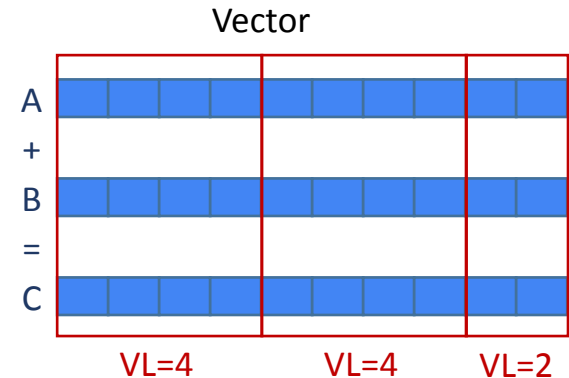




# A bit more elegant: Variable Vector Length

- VL (Vector Length register) is a register (implemented in the hardware) that stores the size of the vector register available at any given point in time
  - VL can be loaded prior to executing the vector instruction with a special instruction (vsetvl)
  - No need to handle “loop tails”
  - Makes the code “vector length agnostic”

```
7 void axpy(double a, double *dx, double *dy, int n) {
8     int i;
9
10    long gvl = __builtin_epi_vsetvl(n, __epi_e64, __epi_m1);
11    __epi_lxf64 v_a = _MM_SET_f64(a, gvl);
12
13    for (i = 0; i < n; i += gvl) {
14        gvl = __builtin_epi_vsetvl(n - i, __epi_e64, __epi_m1);
15        __epi_lxf64 v_dx = _MM_LOAD_f64(&dx[i], gvl);
16        __epi_lxf64 v_dy = _MM_LOAD_f64(&dy[i], gvl);
17        __epi_lxf64 v_res = _MM_MACC_f64(v_dy, v_a, v_dx, gvl);
18        _MM_STORE_f64(&dy[i], v_res, gvl);
19    }
```



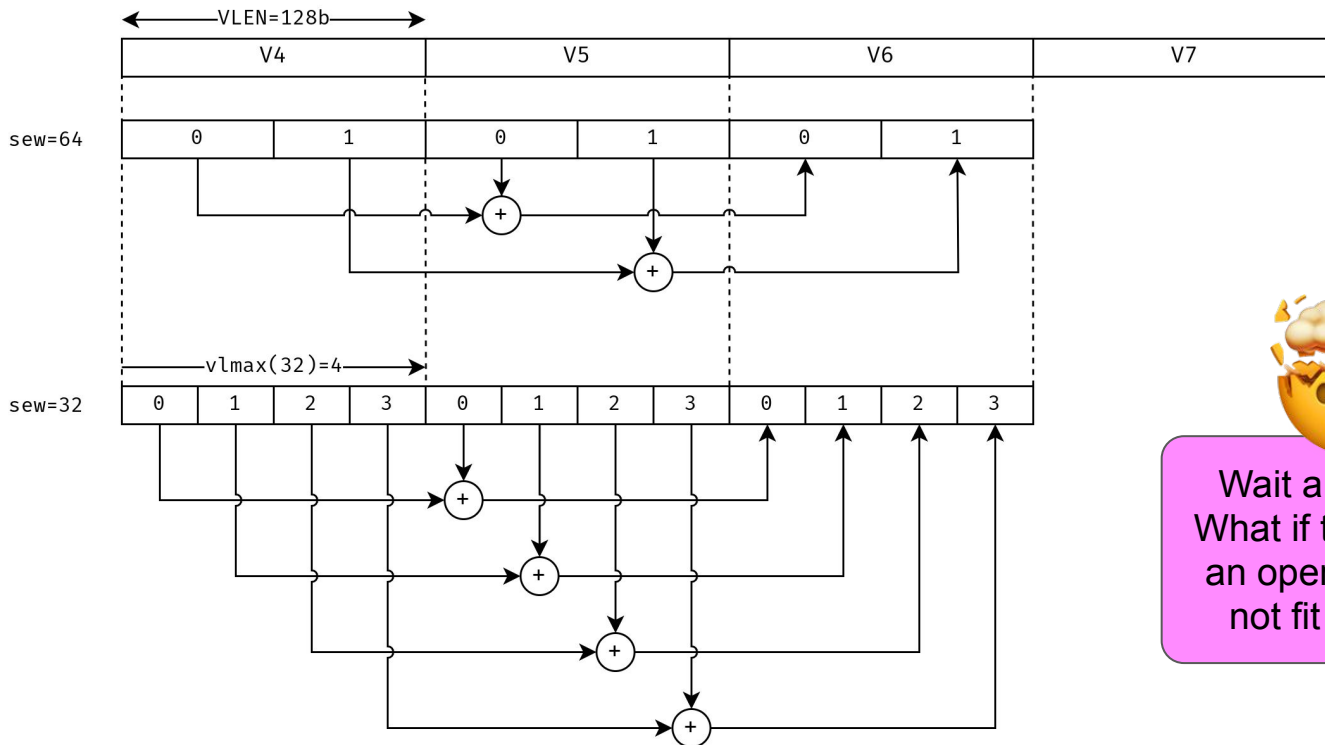
# Take-home message

- RISC-V offers an extension for handling vectors (RVV)
  - It support vectors up to  $2^{16}$  bit register (very large vectors!)
  - VLEN must be defined by the implementor
  - It supports different types of data
  - It supports variable vector length
  - All these features makes RVV Vector Length Agnostic

# An example vector operation with $vl=vlmax$

`vadd.vv v6, v4, v5`

`vl = vlmax(sew)`



Wait a second...  
What if the result of  
an operation does  
not fit the sew?

# Widening and narrowing

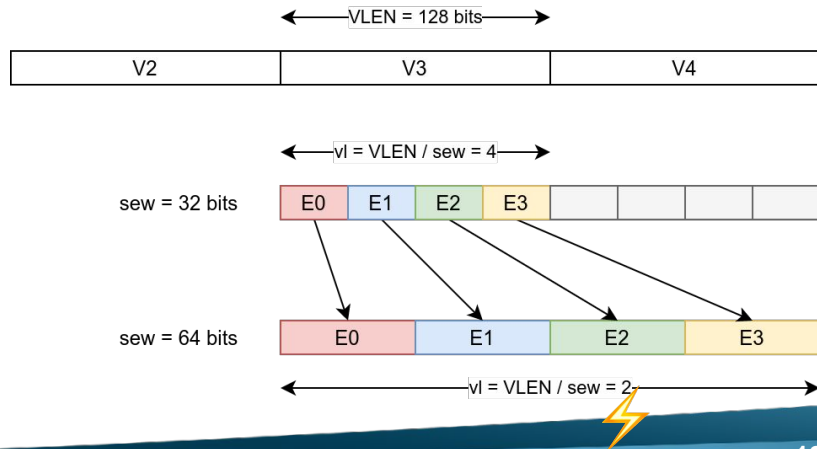
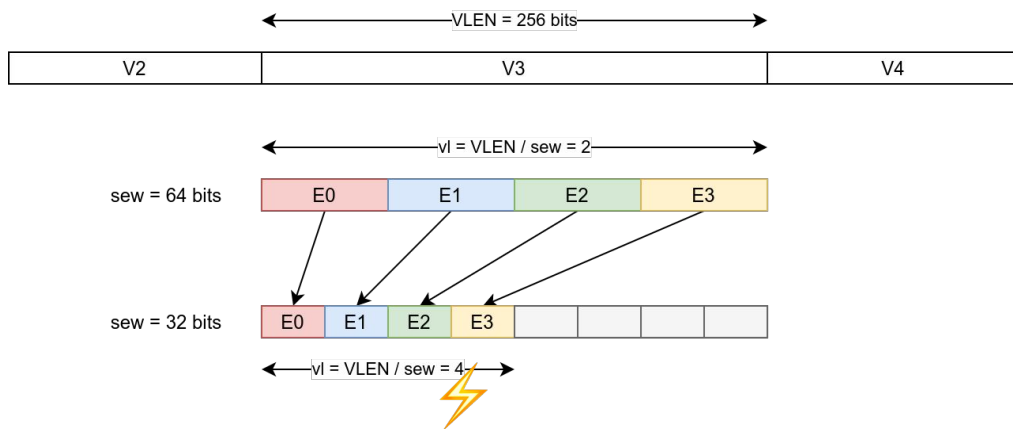
RISC-V RVV provides specific instructions for widening and narrowing operations. Here are some examples:

## 1. Widening Instructions:

- `vwadd.vv`: Vector widening add (e.g., 16-bit → 32-bit).
- `vwmul.vv`: Vector widening multiply (e.g., 16-bit → 32-bit).
- `vfwadd.vv`: Vector widening floating-point add (e.g., 32-bit → 64-bit).

## 2. Narrowing Instructions:

- `vncclip.wv`: Vector narrowing integer clip (e.g., 32-bit → 16-bit).
- `vfncvt.f.x.w`: Vector narrowing floating-point convert (e.g., 64-bit → 32-bit).



# Widening and narrowing: Practical Use Cases

## 1. Machine Learning:

- Widening: Accumulating 16-bit gradients into 32-bit during training.
- Narrowing: Storing 32-bit weights as 16-bit for inference.

## 2. Signal Processing:

- Widening: Performing high-precision filtering on 16-bit audio samples.
- Narrowing: Storing processed audio in 16-bit format.

## 3. Scientific Computing:

- Widening: Avoiding precision loss in intermediate calculations.
- Narrowing: Reducing memory usage for large datasets.

# Widening and narrowing

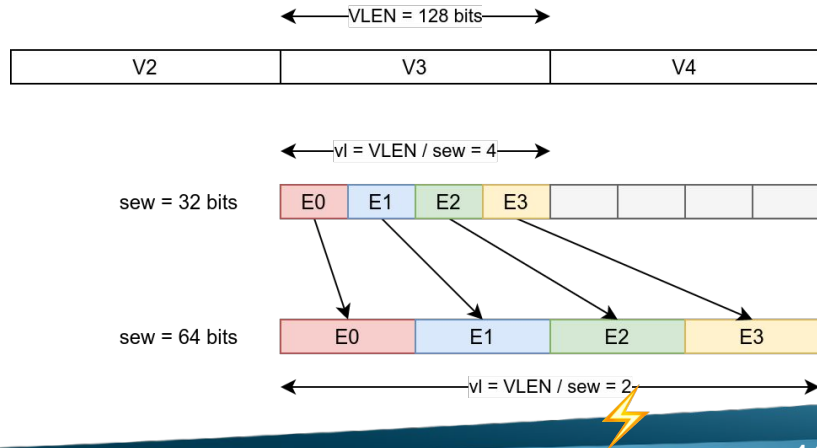
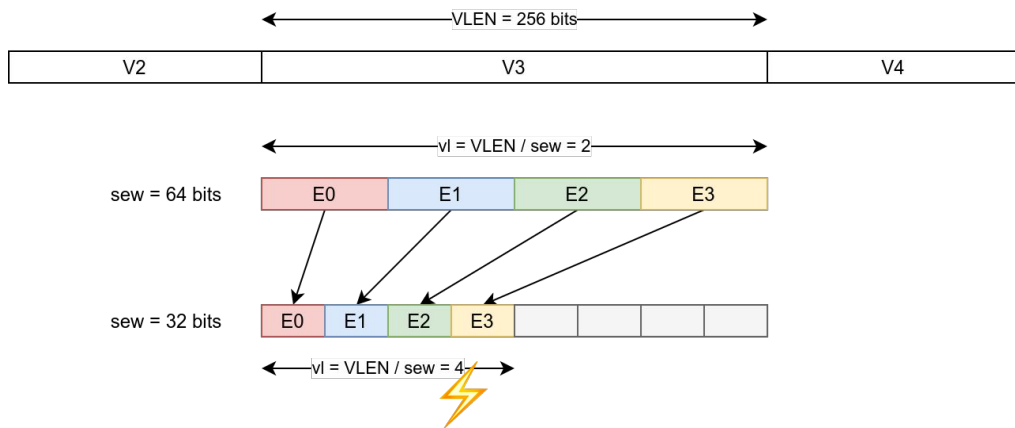
RISC-V RVV provides specific instructions for widening and narrowing operations. Here are some examples:

## 1. Widening Instructions:

- `vwadd.vv`: Vector widening add (e.g., 16-bit → 32-bit).
- `vwmul.vv`: Vector widening multiply (e.g., 16-bit → 32-bit).
- `vfwadd.vv`: Vector widening floating-point add (e.g., 32-bit → 64-bit).

## 2. Narrowing Instructions:

- `vncclip.wv`: Vector narrowing integer clip (e.g., 32-bit → 16-bit).
- `vfncvt.f.x.w`: Vector narrowing floating-point convert (e.g., 64-bit → 32-bit).



# Vector Register Group Multiplier: Imul

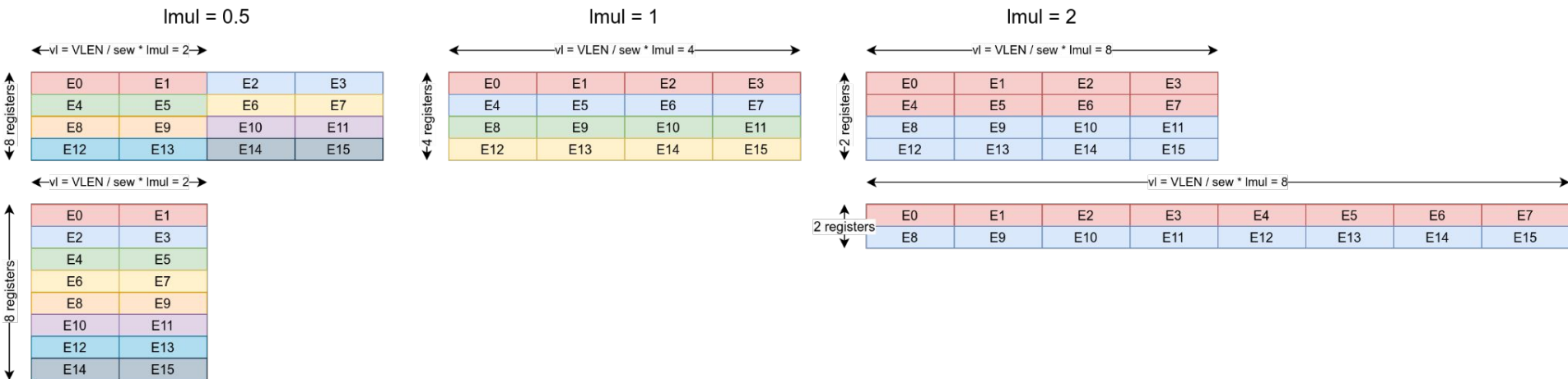
- Imul aka length multiplier
- is a mechanism to cope with different sew
  - if widening makes  $\text{sew} * 2 \rightarrow \text{Imul} = 2$
  - if narrowing makes  $\text{sew} / 2 \rightarrow \text{Imul} = \frac{1}{2}$

Small correction:

- **VLEN**, size of the vectors in bits
- **sew**, Selected Element Width, size of the element of a vector in bits
- **$\text{vmax} = (\text{VLEN}/\text{sew}) * \text{Imul}$** , the number of elements operated by the vector instruction
- **vl**  $\in [0; \text{vmax}-1]$

# Vector Register Group Multiplier: Imul

- RVV allows  $\text{Imul} \in [1/8, 1/4, 1/2, 1, 2, 4, 8]$
- This creates the illusion of
  - a “smaller” register bank with larger registers  $\rightarrow$  when  $\text{Imul} > 1$
  - a “larger” register bank with smaller register  $\rightarrow$  when  $\text{Imul} < 1$





# Take-home message

- RISC-V vector extension allows flexibility with size of the operands
  - This implies re-arrangements using an extra parameter “length multiplier” `lmul`
  - It adds flexibility and the illusion of having more/less registers
  - The full support of `lmul` requires both hardware and software support adding significant complexity to the development of the actual implementation

# Take-home message

- RISC-V vector extension allows flexibility with size of the operands
  - This implies re-arrangements using an extra parameter “length multiplier” `lmul`
  - It adds flexibility and the illusion of having more/less registers
  - The full support of `lmul` requires both hardware and software support adding significant complexity to the development of the actual implementation

**Thank the tech gods somebody took care of the compiler...  
So we don't have to deal with all of this**

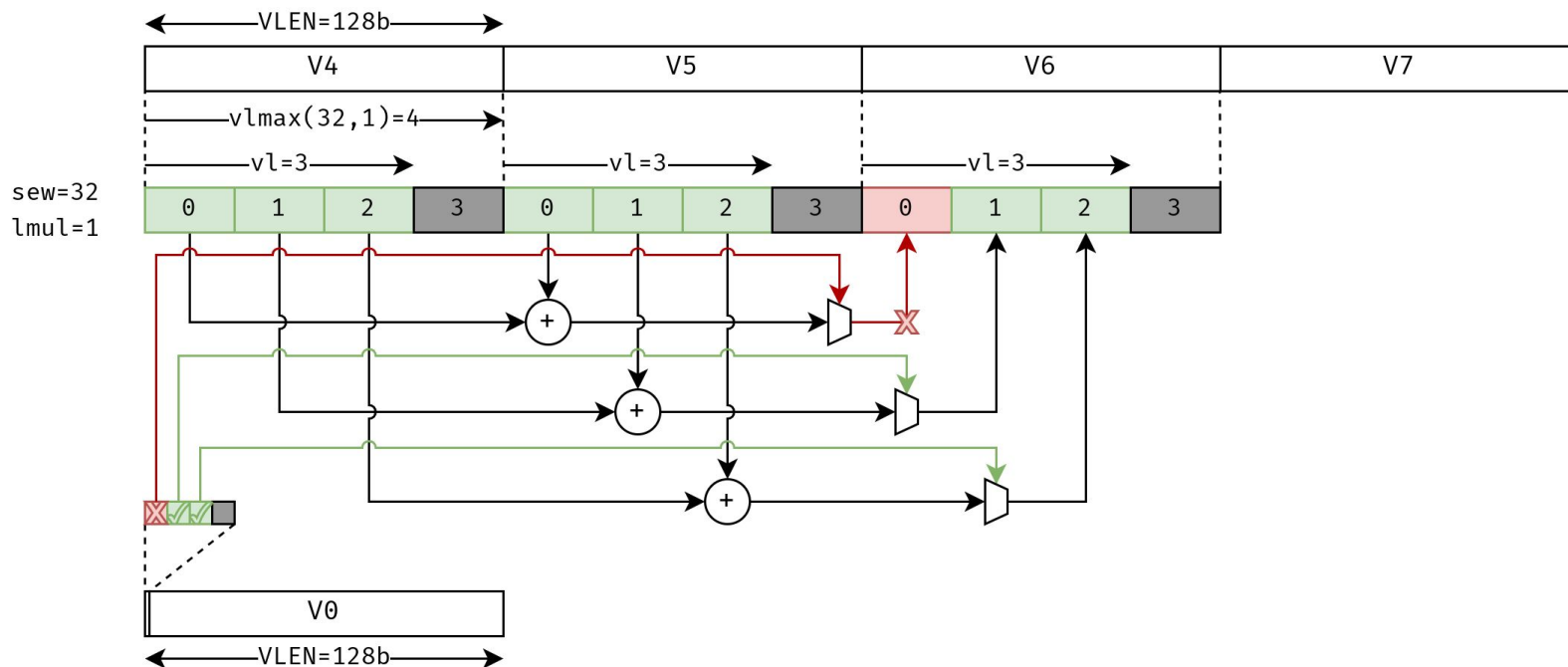


# Masked execution

- A mask vector is a vector of elements of 1 bit in a (regular) vector register (as if  $sew=1$ )
  - An element of the mask set to 0 disables the corresponding element in the vector operation
  - An element of the mask set to 1 enables the corresponding element in the vector operation
- The ISA provides several instructions that generate mask vectors
  - e.g., comparisons, OR-ing two masks, etc.
  - in mnemonics “m” is used to distinguish mask vectors from regular data vectors, marked with “v”
- Instructions can be masked using mask vectors
  - The mask vector of a masked instruction is always v0
- Unmasked instructions can be (conceptually) understood as having a mask of all ones.
  - They may not be equivalent performance-wise though

# Masked execution example

vadd.vv v6, v4, v5, v0.t vl=3



# Take-home message

- RISC-V vector extension improves flexibility allowing the use of masks
  - masks allows to filter elements in a vector operation
  - The full support of masks requires both hardware and software support adding significant complexity to the development of the actual implementation
  - Masks are often generated by conditionals inside loops
    - they introduce performance penalties, but maintain the code linear and clean

**Thank the tech gods somebody took care of the compiler...  
So we don't have to deal with all of this**

