**BFMMLA**

BFloat16 floating-point matrix multiply-accumulate into 2x2 matrix.

If FEAT_EBF16 is not implemented or *FPCR*.EBF is 0, this instruction:

- Performs two unfused sums-of-products within each two pairs of adjacent BFloat16 elements while multiplying the 2x4 matrix of BFloat16 values in the first source vector with the 4x2 matrix of BFloat16 values in the second source vector. The intermediate single-precision products are rounded before they are summed and the intermediate sum is rounded before accumulation into the 2x2 single-precision matrix in the destination vector. This is equivalent to accumulating two 2-way unfused dot products per destination element.
- Uses the non-IEEE 754 Round-to-Odd rounding mode, which forces bit 0 of an inexact result to 1, and rounds an overflow to an appropriately signed Infinity.
- Flushes denormalized inputs and results to zero, as if *FPCR*.{FZ, FIZ} is {1, 1}.
- Disables alternative floating point behaviors, as if *FPCR*.AH is 0.

If FEAT_EBF16 is implemented and *FPCR*.EBF is 1, then this instruction:

- Performs two fused sums-of-products within each two pairs of adjacent BFloat16 elements while multiplying the 2x4 matrix of BFloat16 values in the first source vector with the 4x2 matrix of BFloat16 values in the second source vector. The intermediate single-precision products are not rounded before they are summed, but the intermediate sum is rounded before accumulation into the 2x2 single-precision matrix in the destination vector. This is equivalent to accumulating two 2-way fused dot products per destination element.
- Follows all other floating-point behaviors that apply to single-precision arithmetic, as governed by *FPCR*.RMode, *FPCR*.FZ, *FPCR*.AH, and *FPCR*.FIZ.

Irrespective of FEAT_EBF16 and *FPCR*.EBF, this instruction:

- Does not modify the cumulative *FPSR* exception bits (IDC, IXC, UFC, OFC, DZC, and IOC).
- Disables trapped floating-point exceptions, as if the *FPCR* trap enable bits (IDE, IXE, UFE, OFE, DZE, and IOE) are all zero.
- Generates only the default NaN, as if *FPCR*.DN is 1.

*ID_AA64ISAR1_EL1*.BF16 indicates whether this instruction is supported.

**Vector**
**(FEAT_BF16)**

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 19 18 | 17 16 15 | 14 | 13 | 12 | 11 | 10 | 9 8 7 6 5 | 4 3 2 1 0 |
|----|----|----|----|----|----|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | Rm | 1 1 1 | 0 | 1 | 1 | | | Rn | Rd |

**BFMMLA <Vd>.4S, <Vn>.8H, <Vm>.8H**

```
if !IsFeatureImplemented(FEAT_BF16) then UNDEFINED;
integer n = UInt(Rn);
integer m = UInt(Rm);
integer d = UInt(Rd);
```

**Assembler Symbols**

<Vd>        Is the name of the SIMD&FP third source and destination
            register, encoded in the "Rd" field.

<Vn>        Is the name of the first SIMD&FP source register, encoded
            in the "Rn" field.

<Vm>        Is the name of the second SIMD&FP source register,
            encoded in the "Rm" field.

**Operation**

```
CheckFPAdvSIMDEnabled64();
bits(128) op1 = V[n, 128];
bits(128) op2 = V[m, 128];
bits(128) acc = V[d, 128];

V[d, 128] = BFMatMulAdd(acc, op1, op2);
```

**Operational information**

Arm expects that the BFMMLA instruction will deliver a peak BFloat16
multiply throughput that is at least as high as can be achieved using two
BFDOT (vector) instructions, with a goal that it should have significantly
higher throughput.

---