# arm

# Function Multi Versioning for AArch64:

# Compiler aided function specialization with runtime dispatch

Alexandros Lamprineas

# Intro

## The problem

- Software is deployed on various devices (we may compile for a baseline and reuse the binaries).
- Most CPUs implement optional instructions which may not be present on the target of compilation.
- For example, on Arm, dotproduct instructions may not be available.
- To make use of such instructions, a run-time check is needed.

## The solution

Function Multi Versioning (FMV)
https://learn.arm.com/learning-paths/cross-platform/function-multiversioning

- lets the compiler generate multiple function versions, and
- auto-dispatch between them

# What's new?

- FMV was originally developed for x86 using the function attributes *target* and *target_clones*
  - https://llvm.org/devmtg/2014-10/#talk22      `__attribute__((target("sse4.2")))`
  - https://reviews.llvm.org/D40819

- On AArch64 the *target* attribute is broadly used as an optimization hint,
  - for example in header files of Arm C Language Extensions:

    lib/clang/21/include/arm_acle.h

    ```
    #if defined(__ARM_64BIT_STATE) && __ARM_64BIT_STATE
    typedef struct {
        uint64_t val[8];
    } data512_t;

    static __inline__ data512_t __attribute__((__always_inline__, __nodebug__, target("ls64")))
    __arm_ld64b(const void *__addr) {
      data512_t __value;
      __builtin_arm_ld64b(__addr, __value.val);
      return __value;
    }
    ```

- Therefore, we introduced a new attribute *target_version* (also adopted by RISC-V) and wrote a specification for FMV in ACLE (currently beta).

# Clang CodeGen

- Clang generates multiple function versions with *mangled names* as per
  https://arm-software.github.io/acle/main/acle.html#name-mangling

  - each version may use several features, their names are encoded in the mangled name

  - __attribute((target_version("crc+bti+aes+bf16"))) int fmv(void); → `@fmv._MaesMbf16MbtiMcrc()`

- Each version is associated with *metadata* which propagate information from C/C++ source to LLVM.

- Clang generates a *resolver function* that determines which version to run when the function is called.

  - The resolution is dynamic (it is performed at load time) and permanent for the lifetime of the process.

  - The resolver: (1) initializes the runtime, (2) detects features, and (3) selects the available version of highest priority as indicated by https://arm-software.github.io/acle/main/acle.html#mapping

- Clang generates a global *__aarch64_cpu_features* which contains the feature bits for runtime detection.

arm

# Compiler-rt

The runtime asks the kernel which features are available on host and initializes *__aarch64_cpu_features*

llvm-project / compiler-rt / lib / builtins / cpu_model / aarch64 / fmv / **mrs.inc**

```
1    #if __has_include(<sys/auxv.h>)
2    #include <sys/auxv.h>
3    #define HAVE_SYS_AUXV_H
4    #endif
5
6  ∨ static void __init_cpu_features_constructor(unsigned long hwcap,
7                                                const __ifunc_arg_t *arg) {
8      unsigned long long feat = 0;
9    #define setCPUFeature(F) feat |= 1ULL << F
10   #define getCPUFeature(id, ftr) __asm__("mrs %0, " #id : "=r"(ftr))
11   #define extractBits(val, start, number)
12     (val & ((1ULL << number) - 1ULL) << start) >> start
13     unsigned long hwcap2 = 0;
14     if (hwcap & _IFUNC_ARG_HWCAP)
15       hwcap2 = arg->_hwcap2;
16     if (hwcap & HWCAP_CRC32)
17       setCPUFeature(FEAT_CRC);
103    if (hwcap & HWCAP_SHA3)
104      setCPUFeature(FEAT_SHA3);
105    setCPUFeature(FEAT_INIT);
106
107    __atomic_store_n(&__aarch64_cpu_features.features, feat, __ATOMIC_RELAXED);
108  }
```

## Various platforms are supported

- ∨ 📁 fmv
  - 📄 android.inc
  - 📄 apple.inc
  - 📄 baremetal.inc
  - 📄 freebsd.inc
  - 📄 fuchsia.inc
  - 📄 getauxval.inc
  - 📄 mrs.inc
  - 📄 unimplemented.inc
  - 📄 windows.inc

**arm**

# Example

# Design choices

## Resolver emission

- FMV is supported across multiple translation units.

- The resolver cannot "see" versions beyond the current translation unit.

- Emission options:

  - On use (when the function is called)? → may generate multiple (potentially different) resolvers
    → non deterministic version selection depending on linking order ✖

  - Unique resolver in the TU of the default version
    → deterministic version selection regardless of linking order ✔

## Feature detection

- Dependent-on features get detected transitively as indicated by
  https://arm-software.github.io/acle/main/acle.html#dependencies (for example **sve2**→ **sve** → **fp16** → **fp**)

- Features implied by the command line are not exempt from runtime detection (**simd** → **fp**)

  - For example -march=armv8 implies **simd**

# FMV info representation

FMV info is autogenerated using
llvm/utils/TableGen/ArmTargetDefEmitter.cpp
https://github.com/llvm/llvm-project/pull/113281

llvm-project / llvm / include / llvm / TargetParser / **AArch64CPUFeatures.inc**

```
23      // Function Multi Versioning CPU features.
24  ∨   enum CPUFeatures {
25        FEAT_RNG,
26        FEAT_FLAGM,
27        FEAT_FLAGM2,
28        FEAT_FP16FML,
29        FEAT_DOTPROD,
30        FEAT_SM4,
31        FEAT_RDM,
32        FEAT_LSE,
33        FEAT_FP,
34        FEAT_SIMD,
35        FEAT_CRC,
36        RESERVED_FEAT_SHA1, // previously used and now ABI legacy
37        FEAT_SHA2,
```

- *detection ≠ priority*
  because the detection bit is part of the ABI; if a feature is added/removed whose priority falls between existing ones... 🚩

llvm-project / llvm / include / llvm / TargetParser / **AArch64TargetParser.h**

```
72  ∨   struct FMVInfo {
73        StringRef Name;                  // The target_version/target_clones spelling.
74        CPUFeatures FeatureBit;          // Index of the bit in the FMV feature bitset.
75        FeatPriorities PriorityBit;      // Index of the bit in the FMV priority bitset.
76        std::optional<ArchExtKind> ID;   // The architecture extension to enable.
```

llvm-project / llvm / include / llvm / TargetParser / **AArch64FeatPriorities.inc**

```
16      // Function Multi Versioning feature priorities.
17  ∨   enum FeatPriorities {
18        PRIOR_RNG,
19        PRIOR_FLAGM,
20        PRIOR_FLAGM2,
21        PRIOR_LSE,
22        PRIOR_FP,
```

```
84      // Represents a dependency between two architecture extensions. Later is the
85      // feature which was added to the architecture after Earlier, and expands the
86      // functionality provided by it. If Later is enabled, then Earlier will also be
87      // enabled. If Earlier is disabled, then Later will also be disabled.
88      struct ExtensionDependency {
89        ArchExtKind Earlier;
90        ArchExtKind Later;
91      };
```

- *dependencies*
  are used both for (1) runtime detection, and (2) to enable all the necessary subtarget features for code generation

# Metadata in LLVM IR

- [https://github.com/llvm/llvm-project/pull/118544](https://github.com/llvm/llvm-project/pull/118544)
  Similar to *target-features*.

  **clang/test/CodeGen/AArch64/fmv-features.c**

```
142  + // CHECK: define dso_local i32 @fmv._MaesMbf16MbtiMcrc() #[[unordered_features_with_duplicates:[0-9]+]] {
143  + __attribute__((target_version("crc+bti+bti+bti+aes+aes+bf16"))) int fmv(void) { return 0; }

201  + // CHECK: attributes #[[unordered_features_with_duplicates]] = {{.*}} "fmv-features"="aes,bf16,bti,crc"
```

Why we need them?

- Suppose you have **target_version("i8mm+dotprod")** and **target_version("fcma")**.

- The first version has higher priority because Priority(i8mm) > Priority(fcma) > Priority(dotprod).

- Now suppose you specify **-march=armv8-a+i8mm** on the command line.

- Then the versions would have **target-features** "+dotprod,+i8mm" and "+fcma,+i8mm" respectively.

- If you are using these metadata to deduce version priority, then you would incorrectly deduce that the second version was higher priority than the first!

# GlobalOpt

- May statically (at compile time) resolve calls to versioned functions
  https://github.com/llvm/llvm-project/pull/87939

  by comparing LLVM IR metadata between caller and callee.

*benefit?* → inlining

```
252  + uint64_t AArch64TTIImpl::getFeatureMask(const Function &F) const {
253  +   StringRef AttributeStr =
254  +       isMultiversionedFunction(F) ? "fmv-features" : "target-features";
255  +   StringRef FeatureStr = F.getFnAttribute(AttributeStr).getValueAsString();
256  +   SmallVector<StringRef, 8> Features;
257  +   FeatureStr.split(Features, ",");
258  +   return AArch64::getFMVPriority(Features);
259  + }
260  +
261  + bool AArch64TTIImpl::isMultiversionedFunction(const Function &F) const {
262  +   return F.hasFnAttribute("fmv-features");
263  + }
264  +
```

*metadata selection*

```
58     uint64_t AArch64::getFMVPriority(ArrayRef<StringRef> Features) {
59   +   // Transitively enable the Arch Extensions which correspond to each feature.
60   +   ExtensionSet FeatureBits;
61   +   for (const StringRef Feature : Features) {
62   +     std::optional<FMVInfo> FMV = parseFMVExtension(Feature);
63   +     if (!FMV) {
64   +       if (std::optional<ExtensionInfo> Info = targetFeatureToExtension(Feature))
65   +         FMV = lookupFMVByID(Info->ID);
66   +     }
67   +     if (FMV && FMV->ID)
68   +       FeatureBits.enable(*FMV->ID);
69   +   }
70   +
71   +   // Construct a bitmask for all the transitively enabled Arch Extensions.
72   +   uint64_t PriorityMask = 0;
73   +   for (const FMVInfo &Info : getFMVInfo())
74   +     if (Info.ID && FeatureBits.Enabled.test(*Info.ID))
75   +       PriorityMask |= (1ULL << Info.PriorityBit);
76   +
77   +   return PriorityMask;
78     }
```

*bitmask construction*

- if FMV caller → FMV callee,
  then compare **fmv-features**

- else if non-FMV caller → FMV callee,
  then compare **target-features** with **fmv-features**

# Static resolution algorithm

```
2644  + // Follows the use-def chain of \p V backwards until it finds a Function,
2645  + // in which case it collects in \p Versions. Return true on successful
2646  + // use-def chain traversal, false otherwise.
2647  + static bool collectVersions(TargetTransformInfo &TTI, Value *V,
2648  +                             SmallVectorImpl<Function *> &Versions) {
2649  +   if (auto *F = dyn_cast<Function>(V)) {
2650  +     if (!TTI.isMultiversionedFunction(*F))
2651  +       return false;
2652  +     Versions.push_back(F);
2653  +   } else if (auto *Sel = dyn_cast<SelectInst>(V)) {
2654  +     if (!collectVersions(TTI, Sel->getTrueValue(), Versions))
2655  +       return false;
2656  +     if (!collectVersions(TTI, Sel->getFalseValue(), Versions))
2657  +       return false;
2658  +   } else if (auto *Phi = dyn_cast<PHINode>(V)) {
2659  +     for (unsigned I = 0, E = Phi->getNumIncomingValues(); I != E; ++I)
2660  +       if (!collectVersions(TTI, Phi->getIncomingValue(I), Versions))
2661  +         return false;
2662  +   } else {
2663  +     // Unknown instruction type. Bail.
2664  +     return false;
2665  +   }
```

*discover callee versions*

labrinea marked this conversation as resolved.        ⊕ Show resolved

```
2666  +   return true;
2667  + }
```

```
2686  + static bool OptimizeNonTrivialIFuncs(
2687  +     Module &M, function_ref<TargetTransformInfo &(Function &)> GetTTI) {
2688  +   bool Changed = false;
2689  +
2690  +   // Cache containing the mask constructed from a function's target features.
2691  +   DenseMap<Function *, uint64_t> FeatureMask;
2692  +
2693  +   for (GlobalIFunc &IF : M.ifuncs()) {
```

*for every ifunc in the module*

```
2706  +     // Discover the callee versions.
2707  +     SmallVector<Function *> Callees;
2708  +     if (any_of(*Resolver, [&TTI, &Callees](BasicBlock &BB) {
2709  +           if (auto *Ret = dyn_cast_or_null<ReturnInst>(BB.getTerminator()))
2710  +             if (!collectVersions(TTI, Ret->getReturnValue(), Callees))
2711  +               return true;
2712  +           return false;
2713  +         }))
2714  +       continue;
2715  +
2716  +     assert(!Callees.empty() && "Expecting successful collection of versions");
```

*examine basic blocks of resolver*

```
2725  +     // Sort the callee versions in decreasing priority order.
2726  +     sort(Callees, [&](auto *LHS, auto *RHS) {
2727  +       return FeatureMask[LHS] > FeatureMask[RHS];
```

*sorts callees*

```
2730  +     // Find the callsites and cache the feature mask for each caller.
2731  +     SmallVector<Function *> Callers;
2732  +     DenseMap<Function *, SmallVector<CallBase *>> CallSites;
2733  +     for (User *U : IF.users()) {
2734  +       if (auto *CB = dyn_cast<CallBase>(U)) {
2735  +         if (CB->getCalledOperand() == &IF) {
2736  +           Function *Caller = CB->getFunction();
```

*discover caller versions*

```
2748  +     // Sort the caller versions in decreasing priority order
2749  +     sort(Callers, [&](auto *LHS, auto *RHS) {
2750  +       return FeatureMask[LHS] > FeatureMask[RHS];
```

*sorts callers*

# Static resolution algorithm (continued)

from llvm/test/Transforms/GlobalOpt/resolve-fmv-ifunc.ll

```cpp
2753  +    auto implies = [](uint64_t A, uint64_t B) { return (A & B) == B; };
2754  +
2755  +    // Index to the highest priority candidate.
2756  +    unsigned I = 0;
2757  +    // Now try to redirect calls starting from higher priority callers.
2758  +    for (Function *Caller : Callers) {
2759  +      assert(I < Callees.size() && "Found callers of equal priority");
2760  +
2761  +      Function *Callee = Callees[I];
2762  +      uint64_t CallerBits = FeatureMask[Caller];
2763  +      uint64_t CalleeBits = FeatureMask[Callee];
2764  +
2765  +      // In the case of FMV callers, we know that all higher priority callers
2766  +      // than the current one did not get selected at runtime, which helps
2767  +      // reason about the callees (if they have versions that mandate presence
2768  +      // of the features which we already know are unavailable on this target).
2769  +      if (TTI.isMultiversionedFunction(*Caller)) {
2770  +        // If the feature set of the caller implies the feature set of the
2771  +        // highest priority candidate then it shall be picked. In case of
2772  +        // identical sets advance the candidate index one position.
2773  +        if (CallerBits == CalleeBits)
2774  +          ++I;
```

```cpp
2775  +        else if (!implies(CallerBits, CalleeBits)) {
2776  +          // Keep advancing the candidate index as long as the caller's
2777  +          // features are a subset of the current candidate's.
2778  +          while (implies(CalleeBits, CallerBits)) {
2779  +            if (++I == Callees.size())
2780  +              break;
2781  +            CalleeBits = FeatureMask[Callees[I]];
2782  +          }
2783  +          continue;
2784  +        }
2785  +      } else {
2786  +        // We can't reason much about non-FMV callers. Just pick the highest
2787  +        // priority callee if it matches, otherwise bail.
2788  +        if (I > 0 || !implies(CallerBits, CalleeBits))
2789  +          continue;
2790  +      }
2791  +      auto &Calls = CallSites[Caller];
2792  +      for (CallBase *CS : Calls)
2793  +        CS->setCalledOperand(Callee);
2794  +      Changed = true;
2795  +    }
```

Simplified priority bitmask after dependency expansion:

{mops,sve2,sve,fp16,fp}

- mops+sve2 implies mops → we can statically resolve
- mops implies mops → we can statically resolve bitmask equality → advance callee iterator
- at this point we know the host does not have mops sve does not imply sve2 → we can't statically resolve however sve2 implies sve → advance callee iterator
- keep skipping over callee candidates sve implies sve → advance callee iterator
- no feature is available → we can statically resolve

**sorted caller versions**

→ caller._MmopsMsve2 = {1,1,1,1,1} → callee._Mmops = {1,0,0,0,0} ←
caller._Mmops = {1,0,0,0,0}
caller._Msve = {0,0,1,1,1}
caller.default = {0,0,0,0,0}

**sorted callee versions**

callee._Mmops = {1,0,0,0,0}
callee._Msve2 = {0,1,1,1,1}
callee._Msve = {0,0,1,1,1}
callee.default = {0,0,0,0,0}

arm

# Future work

- User can control feature priorities: https://github.com/ARM-software/acle/pull/371

- User can refer to a specific function version: https://github.com/llvm/llvm-project/issues/84094

- Pointer authentication works with IFUNC resolver: https://github.com/llvm/llvm-project/pull/84704

- Request to support more features (like CSSC): https://github.com/llvm/llvm-project/issues/131218

- ?  (feedback welcome)

**arm**

# Acknowledgements

Special thanks to all the folks who helped with ACLE/code reviews, code refactoring, technical discussions, brainstorming, presentation, etc.

Jon Roelofs, Tomas Matheson, Andrew Carlotti, Daniel Kiss, Wilco Dijkstra, Victor Campos, Sander De Smalen, Maciej Gabka, Andre Vieira, Richard Sandiford, Pavel Iliin, Kristof Beyls, Alfie Richards and others.

arm

Merci
Danke
Gracias
Grazie
谢谢
ありがとう
Asante
Thank You
감사합니다
ধন্যবাদ
Kiitos
شكرًا
ধন্যবাদ
תודה
ధన్యవాదములు
Köszönöm

# arm

arm