

Downloadable Modules Ease Memory Constraints in Small Embedded Systems

Concern over physical size, power consumption and battery life place constraints on memory in deeply embedded systems. These constraints can be eased by using techniques commonly found in large systems but not often seen in small RTOSs.

by John A. Carbone, Express Logic

Embedded systems have always found themselves between a rock and a hard place with regard to the amount of memory available for their software and data. More is always desirable, based on the benefits to functionality, capacity and redundancy, yet at the same time undesirable because of the impact to size, power, heat and cost. Over time, various compromises and advances in technology have changed the picture somewhat, but there has constantly been this tension that will likely always be a limiting factor. After all, Parkinson's Law tells us that memory demand will expand to exhaust available capacity. More memory leads to more demand; hence, we'll always want to stuff more into a system, no matter what the memory limit might be.

Back in the early 1970s, I worked on the Airborne Warning and Control System (AWACS) communications subsystem, an early embedded system based on an IBM 4-pi/ CP-2 processor that used a resounding 8 Kbyte of memory! All our

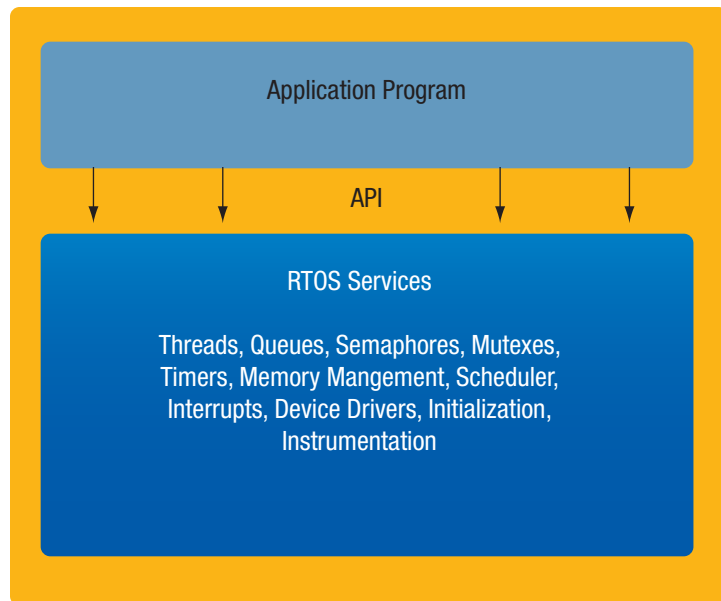


FIGURE 1

Most small RTOSs operate as functions linked with application code into a single executable file.

team's programming was done in assembly language, cross assembled on a DEC PDP-8i, then downloaded into the CP-2's core memory via a monstrosity called a field operating unit (FOU). Since the system had to operate within the 8 Kbyte

limit, there were severe constraints on how the system was designed, what functionality was included, and how it was programmed. We squeezed every instruction out of our code, spending days to trim a few bytes—because we had to.



Get Connected
with companies mentioned in this article.
www.rtc magazine.com/getconnected

We'd have been delirious with even another 8 Kbyte!

As memory constraints have diminished, I often imagine how we would have reacted to the availability of a megabyte. Today, by contrast, we have megabytes in many embedded systems, although some still struggle with less. While today's systems are not quite as limited as 8 Kbyte to be sure, some applications that are particularly sensitive to the same constraints that existed in the 1970s—cost, space and power—may only have perhaps 64 Kbyte to 256 Kbyte. Regardless of the improved memory resources, embedded systems never have enough memory, whether it's due to available technology or size/cost/power constraints.

To meet system functionality requirements, while remaining within the limits of memory availability, embedded developers from time to time borrow big system features to better handle certain situations. Again glancing back to the mainframe days of the 1970s, available memory could be used more efficiently if it were “time shared” among various pieces of independent code. These code sections were called “overlays,” and they were loaded into memory on demand from mass storage, overlaying previously resident code. Of course, this introduced significant overhead, but it was a very effective approach for mainframe systems where such overhead was relatively harmless. With this method, you could introduce virtually unlimited functionality, structured as a set of overlays, without expanding physical memory at all.

A modern technique for functional expansion has worked its way into everyday life through the “app.” An app, or application, is structured in a manner that allows it to be loaded onto a running system. When a system must be updated without removal from service, an app can be downloaded into the system, adding to the memory-resident code that can be run at any time.

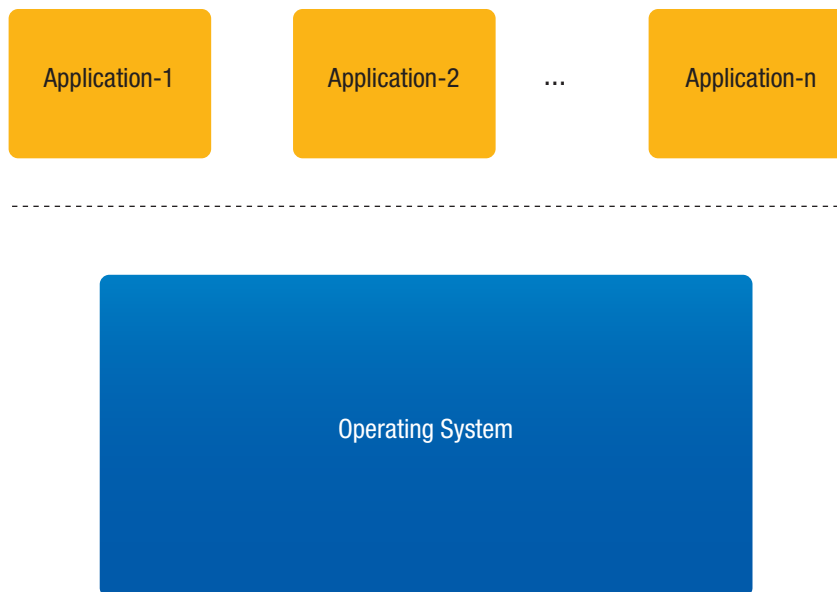


FIGURE 2

A large RTOS, like a desktop OS, is structured as a stand-alone kernel and independent application executables that access it through a trap interface.

Both the overlay and app can be leveraged to give embedded systems a boost in meeting the constraints of a memory system that's just not as big as developers want. In fact, they can be blended in a manner that suits embedded systems quite well, and provides some measure of relief from memory constraints.

Modern Embedded Systems Are Still Memory Constrained

Many of today's embedded systems used in consumer, medical and industrial applications face memory constraints similar to that early AWACS system and many embedded systems since then. These systems run with a real-time operating system (RTOS), and often with relatively small amounts of memory.

RTOSs come in all sizes and flavors, from the large (256 Kbyte - 1 Mbyte), like Wind River's VxWorks, to the super-compact (under 10 Kbyte), like Express Logic's ThreadX. Large RTOSs offer many features adapted from desktop systems that are typically not available in the

super-small RTOSs because such features require a larger amount of code, and result in a slower real-time response. The small RTOS generally operates as a library of services, linked with the application into a single executable file (Figure 1). The application references the services it needs through an API. A small RTOS can provide these services with low overhead and in a small memory footprint, often less than 20 Kbyte. This single executable file is efficient in both time and space since all of the code is statically linked. But it lacks flexibility since any changes to the application or RTOS require rebuilding, relinking, and a redownload/flash of the new image in its entirety.

In contrast, desktop operating systems such as Windows and Linux, and large RTOSs, such as VxWorks, have two-piece “OS/Application” architectures (Figure 2) where target memory is allocated for the OS, and additional memory is allocated for each application that is to be run with the OS, usually as needed.

In this architecture there is a resi-

dent kernel containing all the OS services available to applications, all linked into a distinct, memory-resident executable. This kernel executable boots the system and runs continuously, providing a foundation for applications that are loaded and run in dynamically allocated memory. Of-

gate or to stop the system if it is impossible to perform a requested operation. A divide by zero or the loading of unaligned data illustrate possible illegal operations that might create a trap and activate a trap handler. Similarly, an intentional trap can be generated through an instruction such

lay. This approach produces a time shared memory capability that enables system functionality to exceed physical memory. Such modules can be downloaded into memory via a wireless network or local mass storage, like an app.

The modules use kernel services via an interface with a module manager, an agent within the kernel image that loads and initializes a module as well as fielding all module requests for RTOS services. The module manager is responsible for allocating memory for a module, loading it via a communications or storage interface, and keeping track of its location in memory. The module manager also informs the module of the addresses of the RTOS routines the module's application threads need to access, avoiding the need for a trap mechanism, and providing a very efficient calling interface. The module manager also can overwrite one module with another, if vacant memory is not sufficient to hold the new module—much like the overlays of old. Threads within modules make RTOS service calls exactly as they would make calls if the service function were directly linked with the application. In the module, however, these calls are handled by an interface component that communicates with the module manager. The trap mechanism is avoided, enabling a low overhead service call interface (Figure 5).

Application threads still can be linked with the kernel and reside with the kernel in target memory as part of the RTOS's executable image. This option enables the system designer to balance the need for minimum overhead with memory conservation, using the more efficient, but more memory hungry approach only where absolutely necessary.

Securing Application Modules

Modules also provide a convenient means for memory to be protected against inadvertent, unintended access. They also enable developers to distinguish between more trusted and less trusted modules, applying the appropriate level of protection as demanded by each module. Once a module's code has been proven to be trusted, for example, it can be given ac-

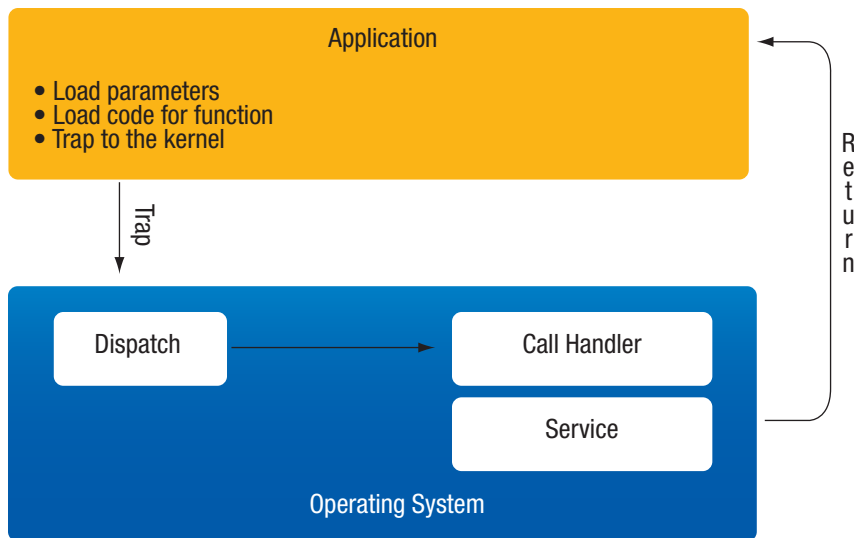


FIGURE 3

A trap mechanism can be used to enable code outside the kernel to communicate with it.

ten in these systems, a virtual memory architecture provides demand paging to and from mass storage on desktop systems or multi-user separation in embedded systems. Virtual memory also can enable use of non-contiguous memory fragments that otherwise might grow over time and reduce or exhaust available memory.

Applications that are not linked with the RTOS must use a different means of accessing RTOS services from the one used by functions within a statically linked executable. This alternate access is necessary because the application is developed independently of a particular hardware environment, and simply cannot predict the memory address of the OS service routine they wish to use. In most cases, these applications access the required OS service routines via a “trap” mechanism that does not require the calling program to be linked with the service it calls (Figure 3).

A trap is a software interrupt that can be used to catch errors before they propa-

as a software interrupt or “swi” instruction, and used to communicate between non-linked program elements.

The trap process requires interrupt servicing, processing, a function call, and then the same in reverse. For a desktop or large RTOS, this overhead is insignificant given the relatively lengthy response time expected of such systems. But for a small RTOS used in hard real-time, deeply embedded applications, low overhead response is essential.

Bringing Dynamic Memory Use to the Small Footprint RTOS

To provide the ability to add functionality dynamically, an “application module” structure can be used (Figure 4). An application module is a collection of one or more application threads, not linked with the kernel, but instead built as a separate executable that can be loaded into target memory when needed, and overwritten when no longer needed—just like an over-

cess to the entire system, including hardware control registers, other threads, and the RTOS, without risk. This trust requires extensive testing and perhaps certification as well. For less critical code, a more practical approach might be to allow that code to access only the memory within its own module and nothing outside.

Despite testing, there is risk of accidental stack or memory corruption due to an erroneously calculated pointer, array limit or stack overflow. These faults can be catastrophic and difficult to find, especially if only one portion of the development team is familiar with the offending module. The extreme difficulty in tracing to the source of such errors makes it all the more important to avoid them. This protects the rest of memory, including other threads and the RTOS, from faults within the un-trusted module code. While all code is tested, a 100% assurance of correctness is expensive and time-consuming. It may be more cost-effective, depending on the code's function in the system, to simply isolate it from other code just in case it misbehaves.

Using the system's MMU or MPU, memory boundary registers can be set to constrain an un-trusted module's code to accesses within a given memory region (Figure 6). With boundary register settings, any attempt to access memory locations outside the specified range will result in a trap. In that event, the trap handler can take appropriate action. It can terminate the offending thread, restart it, alert the system, or halt the system, all depending on the criticality of the offending code. The trap also provides excellent debugging information, as it identifies the code making the errant access, allowing it to be fixed.

Downloadable application modules enable the RTOS to dynamically load and run additional application threads beyond those linked with the kernel. Applications gain increased functionality without the cost of an increased footprint or additional memory, and while retaining an efficient service call interface. This technique also provides on-demand reconfiguration and updates for deployed systems. Downloadable application module technology ideally suits situations where applica-

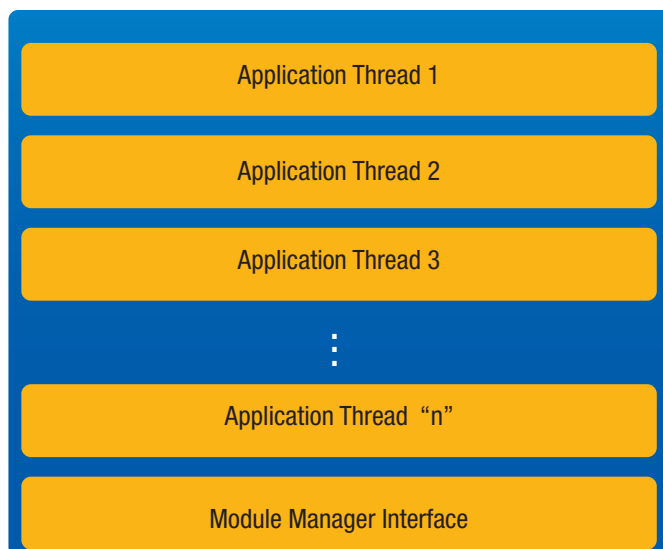


FIGURE 4

Downloadable application modules are executable files containing one to many application threads, with an interface mechanism for accessing kernel services.

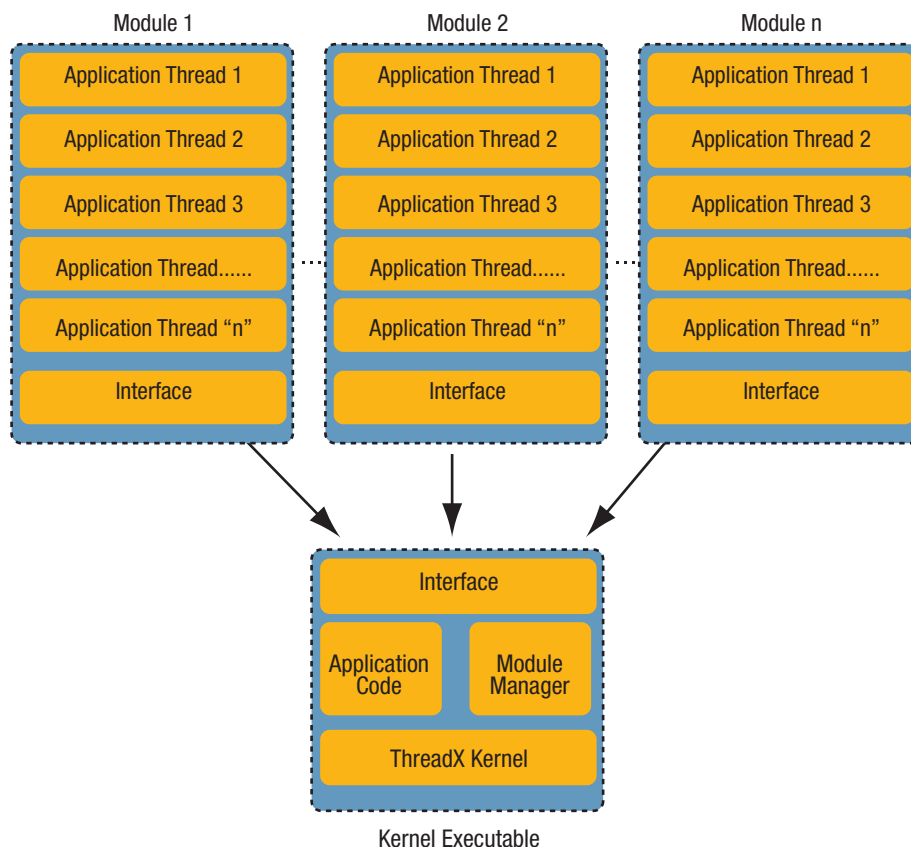


FIGURE 5

Critical applications can be linked with the kernel while others reside in separate modules.

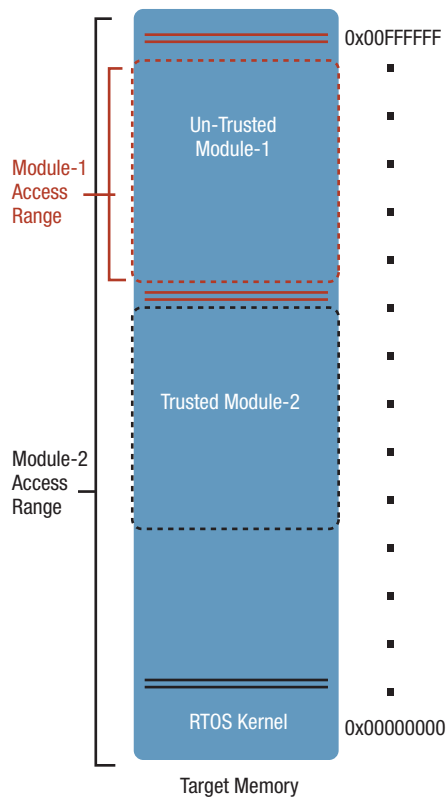


FIGURE 6

Modules and the RTOS can be protected against errant access from other modules, through use of an MMU or MPU.

tion code size exceeds available memory, when new modules need to be added after a product is deployed or when partial firmware updates are required. Another advantage of downloading separate application modules is that each module can be more readily developed by its own team or individual programmer. Each team can then focus on one aspect of a product’s functionality, without having to be concerned with other details. ▲

Express Logic
San Diego, CA.
(858) 613-6640.
www.rtos.com].