

# Asynchronous Systems

Jim Fawcett

CSE681 – Software Modeling & Analysis

Fall 2008

# References

- Concurrent Programming in Java, Second Edition, Doug Lea, Addison-Wesley 2000
  - Several of the slides in this presentation closely follow the material in this book. In some cases the slides simply paraphrase material presented there.
- Programming .Net, Jeff Prosise, Microsoft Press, 2002
  - See Chapters 14 and 15.
- [www.25hoursaday.com/CSharpVsJava.html#threads](http://www.25hoursaday.com/CSharpVsJava.html#threads)
  - Nice discussion of the volatile keyword

# Agenda

- What is an asynchronous system?
- When should you use asynchronous methods?
- Operating system support – Demo code
  - Asynchronous methods
  - Asynchronous delegates
  - Asynchronous callbacks

# Synchronous - Definitions

- **Orbital Satellites**
  - Stationary relative to the earth
- **Neurobiology**
  - Mental processes that entrain to external stimuli as in epilepsy
- **Communication Systems**
  - Information is contained in frames with constant frame rate
- **Radio and Radar Detection**
  - Carrier is removed by an oscillator that locks onto the incoming carrier frequency
- **Software**
  - A function call that blocks the caller until finished
  - A component that collects input by waiting for data from a single sender at some point in the code, e.g., cin

# Asynchronous Software

- Function call returns immediately without waiting for function completion.
- One of the following happens after return:
  - Caller needs no reply and ignores the callee.
  - Caller must poll for completion status, e.g., keep checking
  - Caller must provide a callback for the callee to use when finished
  - Caller deposits a message in a queue for the callee to process at some later time, without expecting or waiting for a reply.
    - The callee may, but is not required to, deposit a reply in a queue owned by the caller.
  - If a system is based on message passing, the callee can react to inputs from an arbitrary number of sources, arriving in any order, at any time.
    - Very flexible.

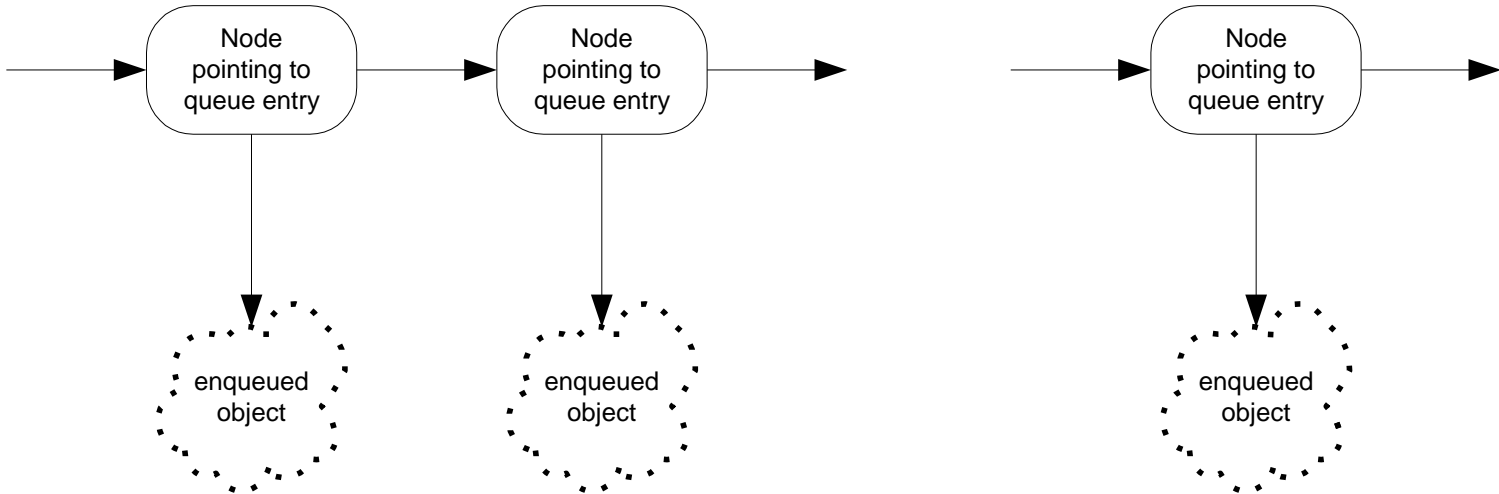
- # Synchronize (synonym – serialize)
- Just to be confusing, the word synchronize, pronounced and spelled very like synchronous, means something entirely different.
  - For software, synchronize means to control access to a resource shared between threads so that only one thread at a time is allowed to use the resource.
    - The resource is locked by the thread and no other thread gets access, even if the thread is blocked, suspended, or waiting to run, until the using thread releases the lock.
    - A lock may apply only to some specific code location or to an object accessed anywhere in the code.
    - A thread can lock a resource by using any of the following:
      - C# lock, .Net Monitor, critical section, mutex, or event
    - The following resources are often shared between threads:
      - Queues, I/O streams, files, and windows
      - Static members of a class
      - C and C++ can share global variables

# What is an Asynchronous System?

- Two parties communicate without being bound to a specific time.
  - Email is a classic example.
    - Message is sent, can be read at any later time.
    - No constraints on when message is sent.
    - No constraints on when message is read, as long as it is later than sending time.
  - Requires four things:
    - Sender
    - Receiver
    - Place to put messages
    - Transmission facility that does not require any action on part of sender or receiver, other than to send and collect message.

# FIFO Queues

- First In First Out queues are usually constructed with linked lists.
  - Objects enter the queue by getting linked to one end.
  - Objects leave the queue by getting unlinked from the other end.

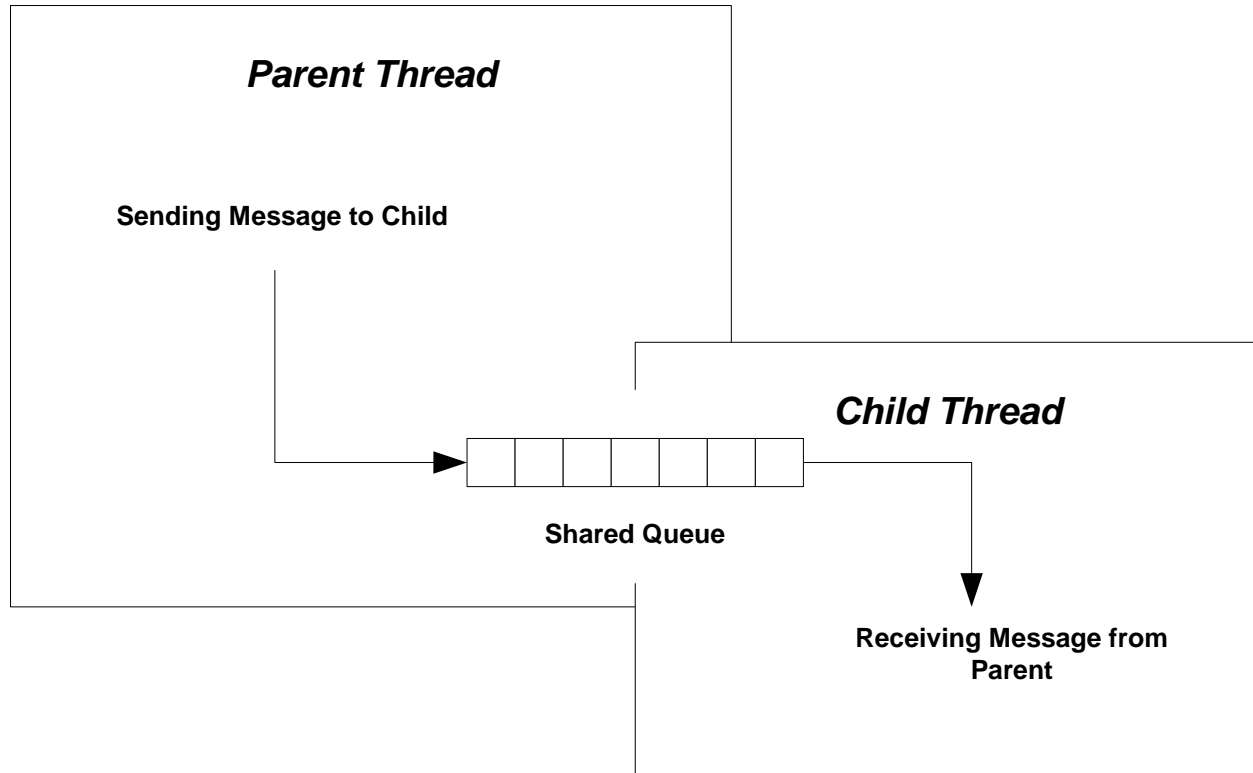




# Important Property of Queues

- Queues decouple a receiver from its sender.
  - Sender and receiver can be on different threads of a given process.
  - Receiver does not need to process an object when it is handed off by the sender.
  - Queues can eliminate timing mismatches between sender and receiver.
    - One might be synchronized in nature, requiring message passing at fixed time intervals – a radar signal processor for example.
    - The other might be free-running, handling messages with different service times, preventing synchronized operation.
- Queues can support reliable communication over unreliable media, simply holding onto messages it can't send until the transmission link becomes available.

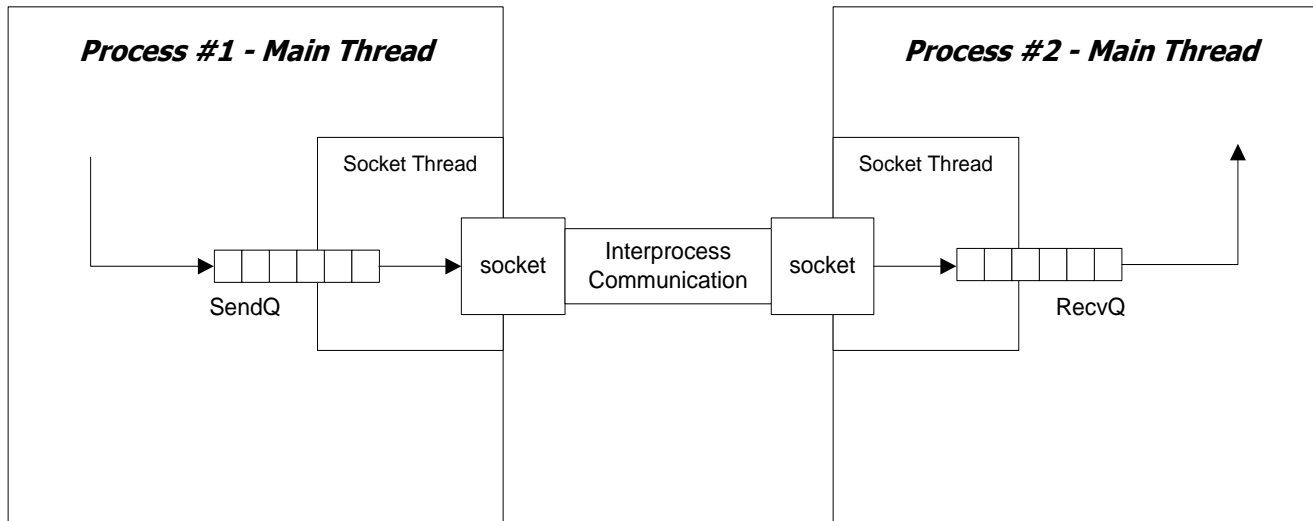
# Message Passing Between Threads



# Send and Receive Queues

- Essentially, a SendQ lets:
  - Sender thread create a list of requests.
  - Communication thread remember requests it has not processed yet.
- A RecvQ lets:
  - Remote process return a result independently of the receiver's main thread's readiness to accept it.
  - Valuable remote resource need not block waiting for a hand-off to receiver.
- Both queues allow the client's main thread to service other important tasks as well as interact with the remote resource.

# Send and Receive Queues



SendQ is used to hold messages in the event that communication with the remote receiver fails.

Messages are held until communication is re-established.

RecvQ is used to quickly remove messages from the socket connection so that the socket buffer never fills (that would block the sender).

# What is an Asynchronous System?

- So, is every message-passing system asynchronous? No.
  - Exchange between a browser and web server is message-based, usually employing Get or Post HTTP messages.
  - That exchange is synchronous. The browser doesn't return until a page is delivered – either the page requested or an error page.
- Is every asynchronous system message-based? No.
  - .Net delegates and remoting proxies support asynchronous operations via BeginInvoke and EndInvoke procedure calls.
- Examples of asynchronous systems:
  - Windows operating system supports the ability to react to many kinds of events using each window's message loop.
  - Socket Listeners improve their availability by spawning client handler threads.
  - All the radar systems I worked on use asynchronous messaging between layers.
  - Email, Project #5, many enterprise systems, ...

# Why Use Asynchronous Systems?

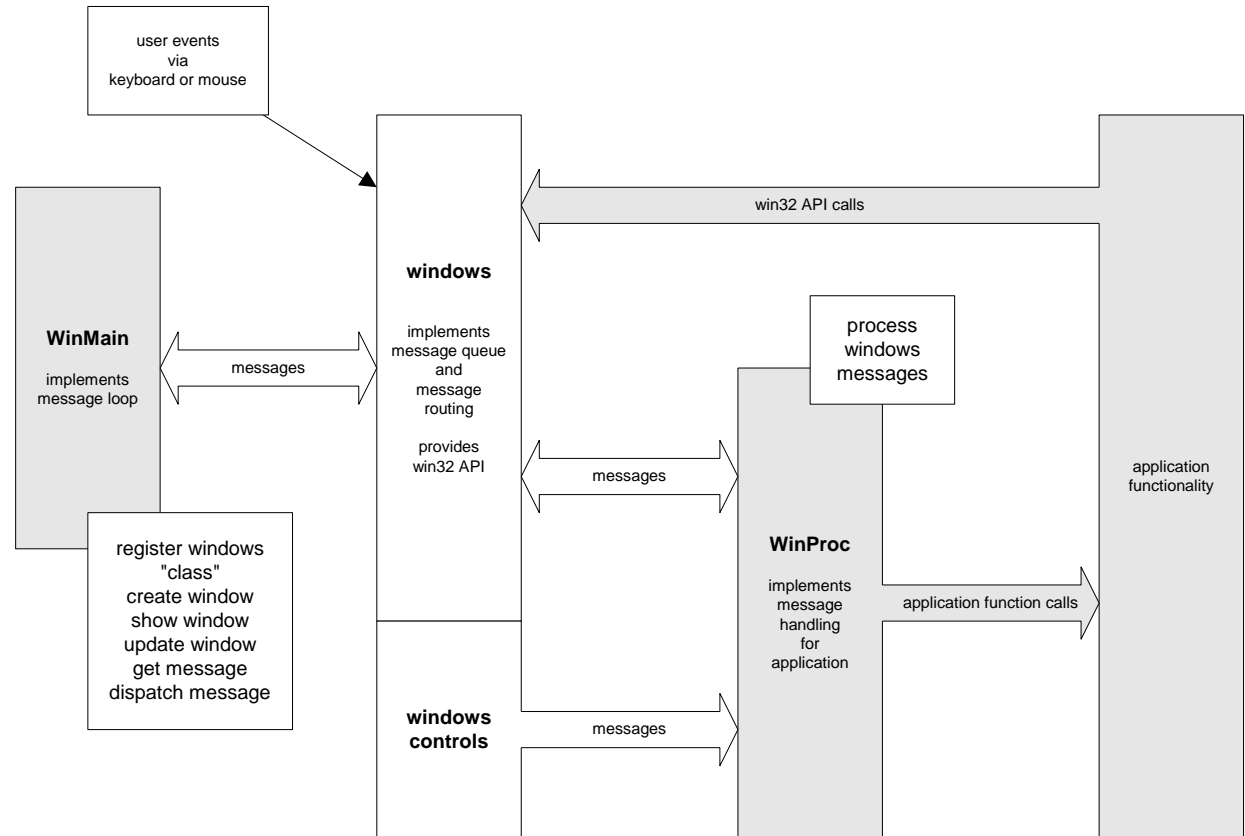
Adapted from Concurrent Programming in Java, Doug Lea, Addison-Wesley, 1997

- System needs to be reactive
- System must have high availability
- Services must be controllable
- System needs to send asynchronous messages
- System may have to handle bursty events

# Why use Asynchronous Systems?

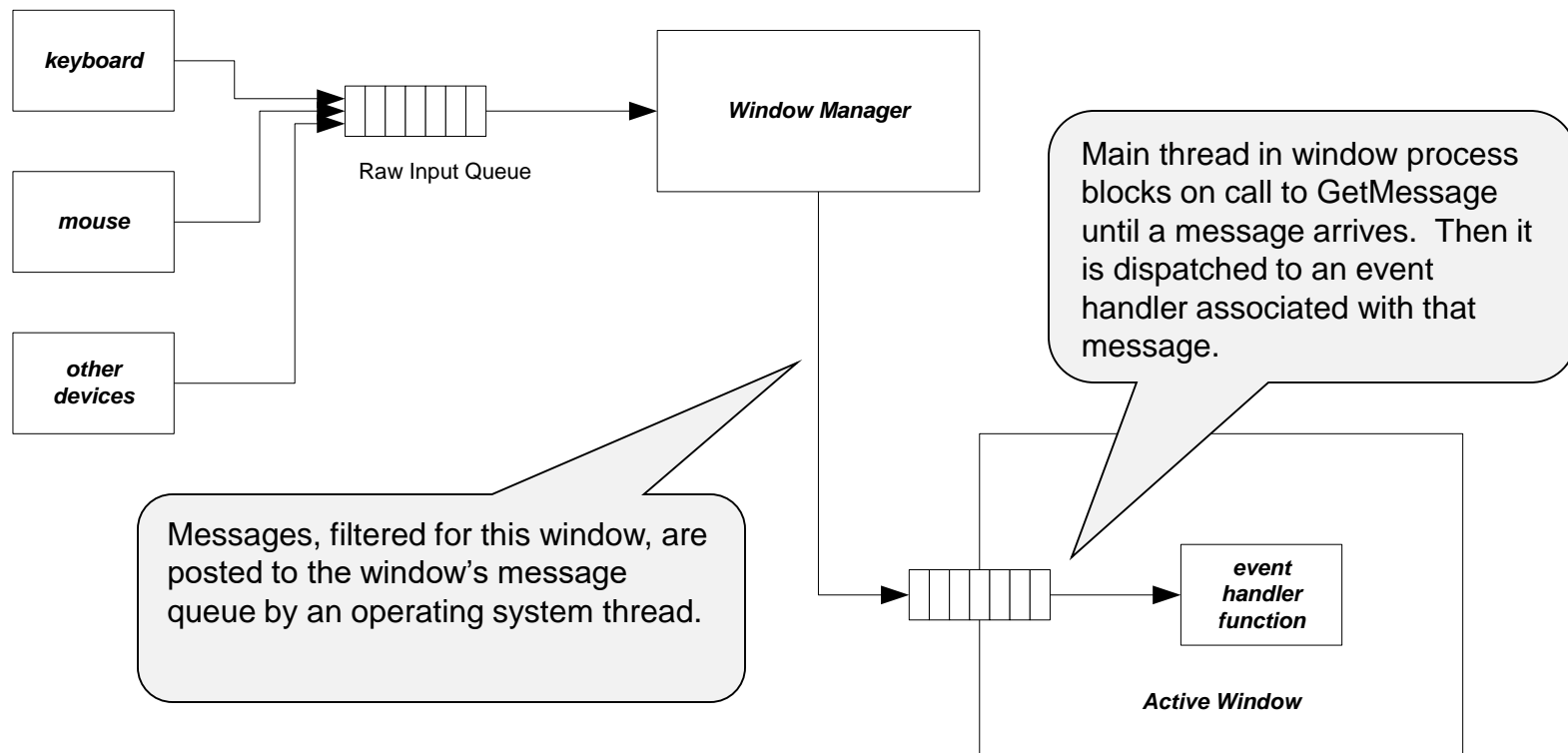
- The system needs to be **reactive**.
  - It does more than one thing at a time, each activity reacting in response to some input.
  - Each event results in a message in windows queue. Reponse happens later, if at all.

## *Classic Windows Processing*



# Windows, Queues, and Messages

- Graphical User Interfaces are the stereotype of message-passing systems using queues.



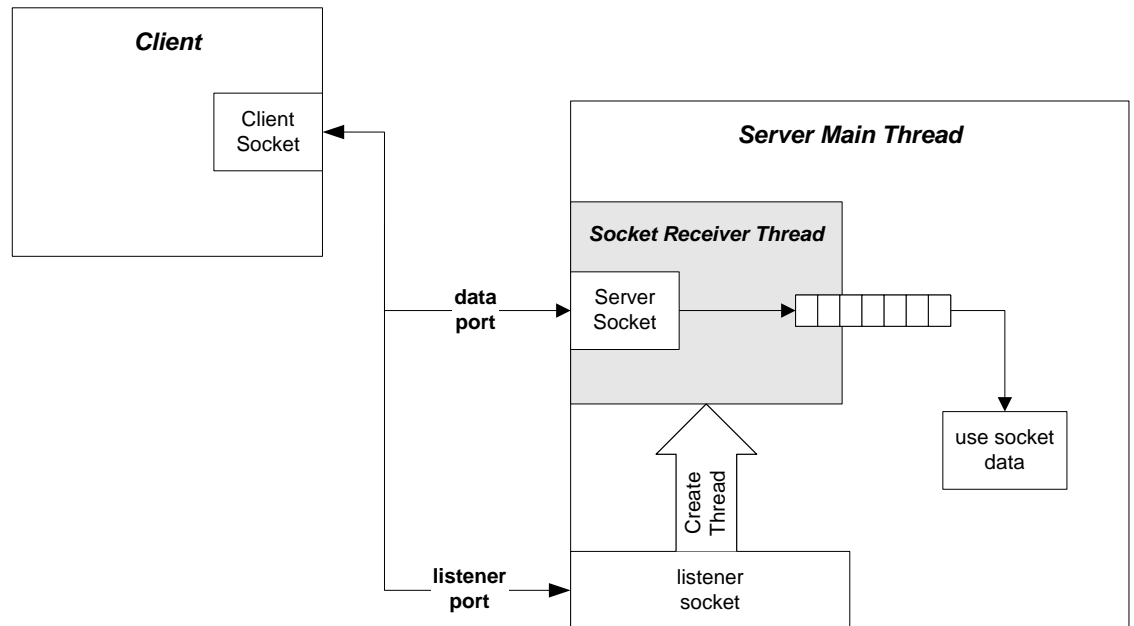


# Windows Messaging

- With the architecture shown on the previous slide a window can respond to any of many different inputs, including:
  - User inputs from mouse movement, mouse buttons, and keyboard
  - System events, e.g., being uncovered by an obscuring window, and so needing to repaint its region of the screen
  - Message generated from within the program running in that process, based on strategies put in place by the designer.
- Even if several messages arrive before a predecessor is processed, they won't be lost.

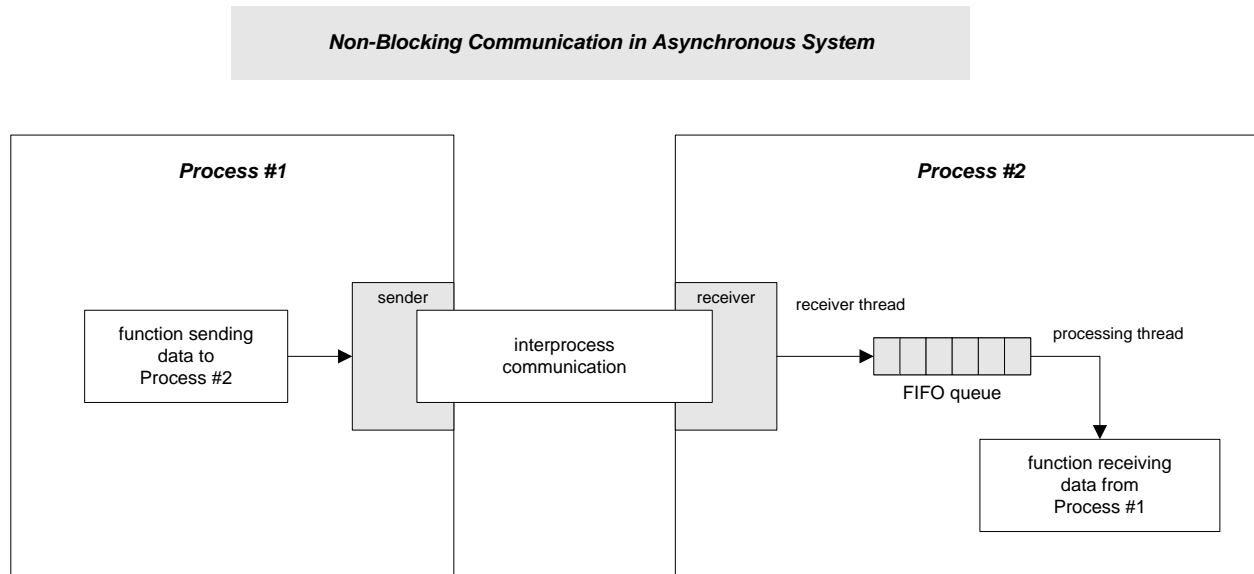
# Why use Asynchronous Systems?

- The system must have **high availability**.
  - One object may serve as a gateway interface to a service, handling each request by constructing a thread to asynchronously provide the service.
  - Listener must quickly dispose of the connection so that it can go back to listening for other requests to connect and clients find the server available.



# Why use Asynchronous Systems?

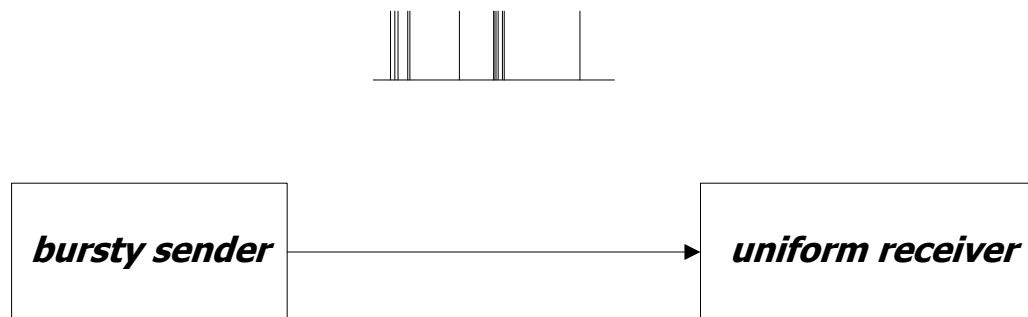
- The services need to be **controllable**.
  - Activities within threads can be suspended, resumed, and stopped by other threads. Controller issues command but does not wait for action.
- The system needs to **send asynchronous messages**.
  - The calling object may not care when a requested action is performed, at least within some reasonable bounds.



# Why use Asynchronous Systems?

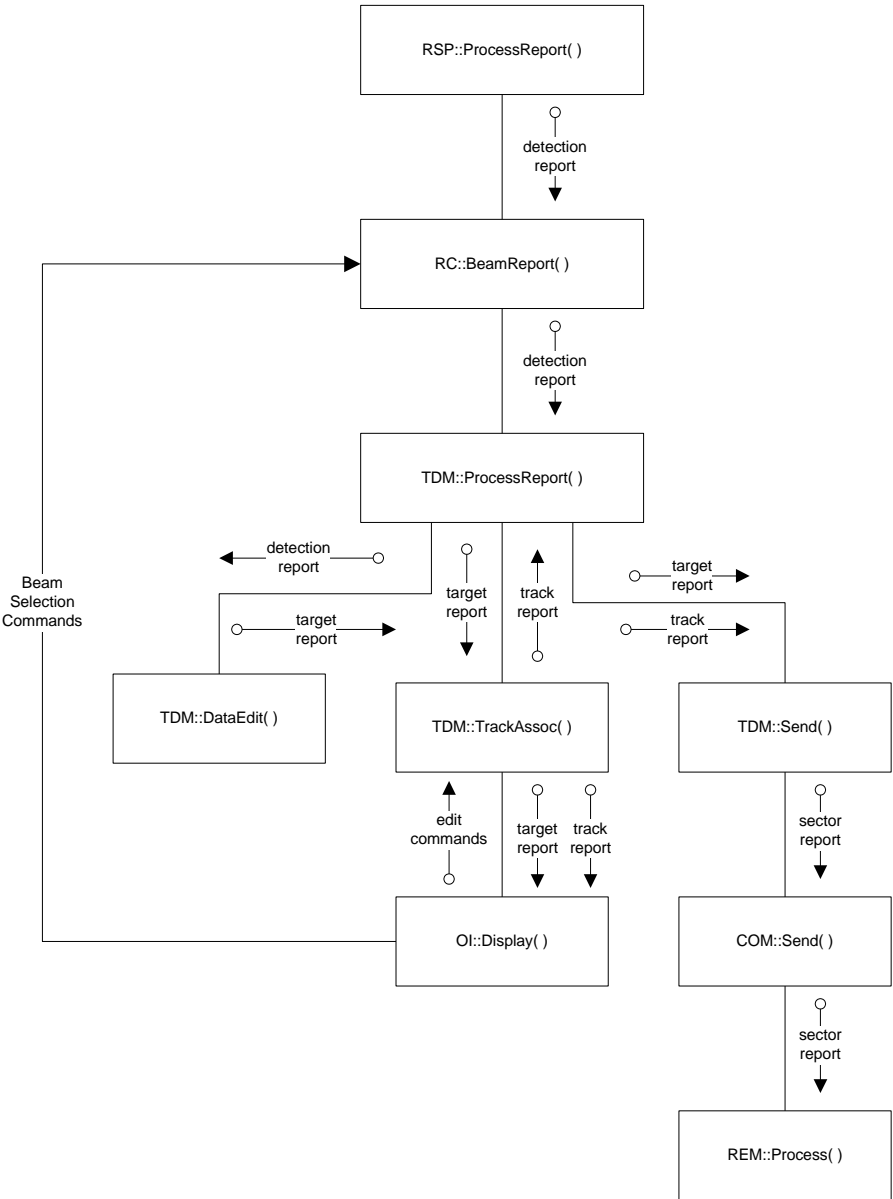
- The system may need to **handle bursty events**.
  - Some events may occur in rapid succession for brief periods, perhaps far faster than the system can react to them. When this happens an asynchronous system can enqueue requests to work on at a later time.

## ***Bursty System***



# Synchronous Radar Just won't work!

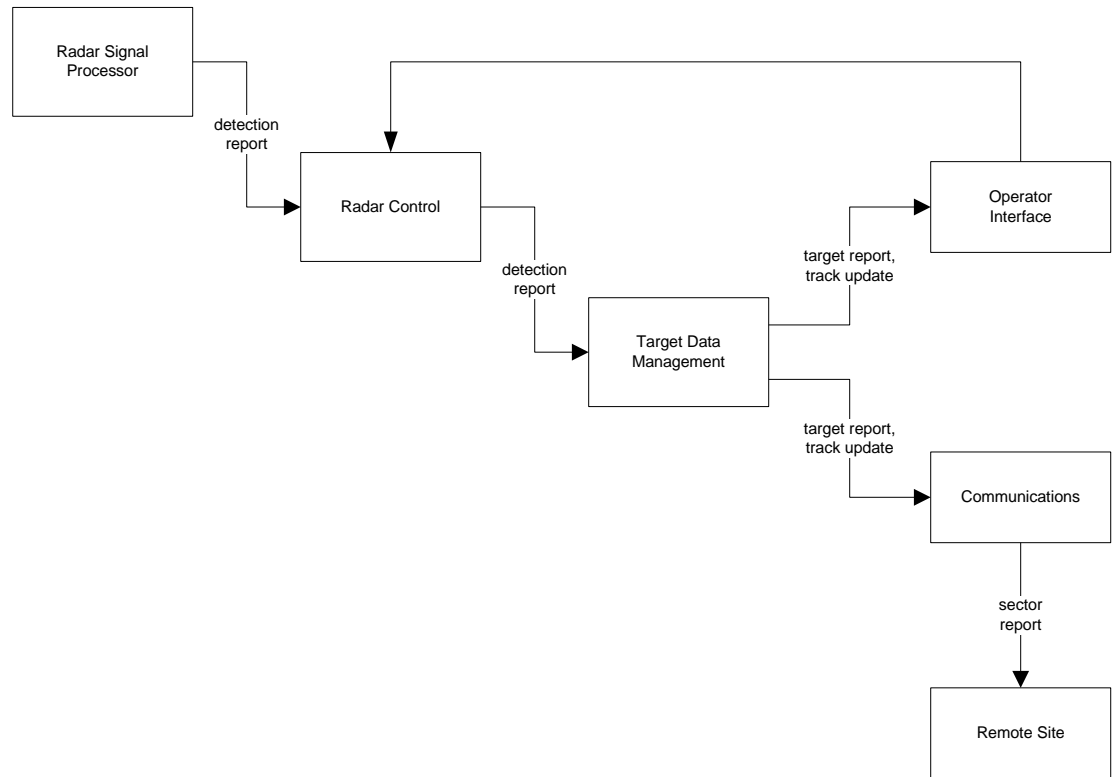
- A synchronous radar would require the signal processor to wait, when it sends a report, for the report to be processed at several levels.
- This just will not work. The signal processor must operate at a high rate of speed to service all of the cells in coverage.
- Each radar component must function independently on the inputs provided to it.



# Asynchronous Radar

- In the asynchronous radar each component computes its output and deposits it in a queue for processing by the next layer.
- The messages at the top of the chain are small and arrive at a high rate.
- Message lower in the chain are larger and less frequent.

## Typical Distributed System - Radar Processing



# Asynchronous Operation

- There are several ways to support Asynchronous Operation. All of these use asynchronous methods:
  - Create a new thread
    - Pass it the method to run asynchronously, using a ThreadStart delegate.
  - Use a thread from a ThreadPool
    - Queue a work item with a delegate pointing to the method to run asynchronously.
  - Use a BackgroundWorker in a Form
    - Has useful events for signaling progress and completion.
  - Use a delegate's BeginInvoke(), EndInvoke() methods
    - Runs the delegate's methods asynchronously.
  - Use a Form's BeginInvoke(), EndInvoke() methods.
    - Allows a worker thread to call a Form method running on the Form's UI thread.
  - Use a class that implements ISynchronizeInvoke interface
    - The ISynchronize interface declares methods:  
InvokeRequired(), BeginInvoke(), EndInvoke(), Invoke()

# Asynchronous Delegate

- `class Asynch`

```
{  
    public delegate string SlowCallDelegateType(int secs);  
    public SlowCallDelegateType callDelegate;  
  
    public string SlowCall(int millisecs) { ... }  
    // more class functionality  
}
```

- `asOp.callDelegate =`

```
    new Asynch.SlowCallDelegateType(asOp.SlowCall);
```

- `// begin asynchronous operation`

```
IAsyncResult ar  
    = asOp.callDelegate.BeginInvoke(1000,null,null);
```

```
    // do some useful work here, then wait  
    // on asynchronous call to finish
```

```
string result2 = asOp.callDelegate.EndInvoke(ar);
```

- Here, `asOp` is an instance of `Asynch`.



# Delegate's BeginInvoke Arguments

- The arguments of BeginInvoke depend on the signature defined by the delegate.
  - The first arguments are the signature arguments ordered as:
    - in parameters
    - out parameters
    - in/out parameters
    - ref parameters
  - These are followed by the final two arguments of type:
    - AsyncCallback
    - AsyncState
- The return value of BeginInvoke is an instance of IAsyncResult.

# Delegate's EndInvoke Arguments

- The arguments of EndInvoke consist of the non-in parameters and the IAsyncResult instance returned by BeginInvoke:
  - The first arguments are the signature arguments ordered as:
    - in parameters *omitted*
    - out parameters
    - in/out parameters
    - ref parameters
  - These are followed by the final argument of type:
    - IAsyncResult instance returned by BeginInvoke.
- The return value of BeginInvoke is the return value of the delegate's function.

# Design Forces

- The primary design forces are safety and liveness:
  - Safety
    - nothing bad ever happens to an object or resource
  - Liveness
    - Something eventually happens within an activity
  - Performance
    - How soon and quickly are services provided?
  - Reusability
    - How easy is it to use services in another programming context?

# Safety

- Read/Write conflicts
  - One thread attempts to read a field while another writes to it.
- Write/Write conflicts
  - Two threads attempt to write to the same field.

# Liveness

- In a live system we expect to make progress toward completion. This may not happen immediately for the following reasons:
  - Locking
    - A lock or synchronized method blocks one thread because another thread holds the lock.
  - Waiting
    - A thread blocks, on a `join()` or `wait()` method, waiting for an event, message, or result computed by another thread.
  - Input
    - An IO method waits for input that has not arrived yet.
  - Contention
    - A runnable thread fails to run because other threads are occupying the CPU or other needed resources.
  - Failure
    - A method encounters an exception or other fault.

# Liveness (continued)

- A permanent lack of progress may result from:
  - Deadlock
    - Circular dependencies between locks
  - Missed event
    - A thread starts waiting for an event after the event occurred
  - Nested lockout
    - A blocked thread holds a lock needed by another thread attempting to wake it up.
  - Livelock
    - A continuously retried action fails continuously.
  - Starvation
    - Scheduler fails ever to allocate CPU time to a waiting thread.
  - Resource exhaustion
    - A thread attempts to access one of a finite set of resources, all of which are in use (file handles for example).
  - Distributed system failure
    - A remote machine, hosting a needed service, is unavailable.

# Performance

- Performance is usually described by the metrics:
  - Throughput
    - Number of operations per unit time, e.g., messages processed, files sent, ...
  - Latency
    - Elapsed time between a request for service and its satisfaction.
  - Capacity
    - The number of simultaneous activities that can be supported for a given throughput or maximum latency.
  - Availability
    - Number of simultaneous requesters that can be supported without failures to connect.
  - Efficiency
    - Throughput divided by the amount of CPU resources needed, e.g., CPUs, memory, IO devices, ...
  - Scalability
    - Rate at which throughput or latency improves when resources are added to the system.
  - Degradation
    - Rate at which latency or throughput decreases as more clients or activities are added without adding new resources

# Performance (continued)

- Factors that affect performance
  - Partitioning of distributed resources
    - Affects the amount and sizes of data sent between processes, machines, and networks.
  - Caching
    - Holds data that has been sent previously for possible use later. Attempts to avoid retrieving resources already held. May induce problems with consistency and staleness of the data.
  - Locking strategy
    - May trade off the use of immutable objects against use of mutable locked objects, e.g., create new immutable objects to achieve a state change, rather than locking and modifying a mutable object's state.
  - Use of background tasks
    - Define activities that can usefully proceed while main activities are blocked.
  - Use algorithms specifically designed for concurrency
    - Some efficient sequential algorithms do not lend themselves to efficient concurrent implementation, but there may be good concurrent versions available.



# Threading Memory Model

- Both the compiler and the processor that your code runs on are allowed to:
  - Cache variables in registers or other cache memory
  - Rearrange instructions
  - For single threads, the compiler and processor are constrained to preserve sequential semantics, e.g., program order semantics.
  - For interactions between threads there are no such guarantees.
- The ***lock*** construct and ***volatile*** qualifier are intended to extend these guarantees to threads operating in a multi-threaded environment.
  - Any thread running within a lock body or in a synchronized function will enjoy the program order semantics guarantee.
    - Releasing a lock forces a flush of all writes from the thread's working memory.
    - Acquiring a lock forces a reload of all fields accessible to the acquiring thread.
  - If a variable is qualified as volatile:
    - Any written value is flushed by the writer thread before the writer performs any other memory operation.
    - Reader threads must reload the values of volatile variables for each access.

End of Presentation

# Three Threaded Client with SC/PBRO

