Multithreaded applications (1)

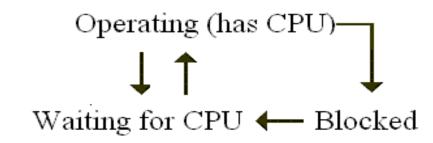
A thread of execution is the smallest sequence of programmed instructions that an operating system can manage independently.

Multitasking: several processes run concurrently (or seem to run concurrently). Between processes the computer resources (like memory, ports, etc.) are not shared, For example each process has its own address space.

Multithreading: several threads run concurrently (or seem to run concurrently) within one process. Several threads can share the same resource, for example a section of memory.

On single processor systems the processor switches from one thread to another. When the time slice allocated for a thread has expired, the processor starts to execute instructions from another thread. As the time slices are short, the user percieves that the threads are running on the same time. On multiprocessor (multicore) systems each processor may run a particular thread, i.e. the threads are actually running on the same time.

A thread may be in three states:



When a thread needs a resource currently occupied by another thread or waits for data from external sources, it is in the blocked state. The other threads continue to work.

Multithreaded applications (2)

When you must create threads:

- 1. You have some subtasks that take a large amount of time (searching, printing, etc.).
- 2. You need to communicate with external data sources (web, devices connected to your computer, etc.) which may send data to you or read data you want to send them at occasional moments. It may even happen that they do not want to communicate with you. A single-thread application waiting for data from external source freezes and becomes uncontrolable.
- 3. You have subtasks with different priority.
- 4. You have user interface which must be responsive i.e. react to the user actions (for example mouse clicks) immediately.

main() is in the primary (or main) thread. The main thread may initialize and launch another threads, those in turn can also have their own child threads and so on.

Each thread must have its thread entry point function which corresponds to the *main()* of the primary thread.

Each thread has its own stack for local variables. The global variables are common to all the threads.

Threads (1)

```
#include <thread> // See <a href="https://en.cppreference.com/w/cpp/thread/thread.html">https://en.cppreference.com/w/cpp/thread/thread.html</a>
using namespace std;
To declare a thread and launch it, write:
thread thread object name(entry_point_function_name,
                              list of input parameters);
The entry point function may have any set of input parameters, but no return value.
Example: suppose we have function
void PrintTextFile(string FileName);
To organize the printing in background write
thread BackgroundPrinting(PrintTextFile, string("c:\\temp\\data.txt"));
The thread stops when its entry point function quits. The question is what happens if the
thread object is destroyed (automatically as a local variable or by delete operator) before
the end of entry point function. For example:
int main()
  thread BackgroundPrinting(PrintTextFile, string("c:\\temp\\data.txt"));
  return 0; // still printing, but BackgroundPrinting thread object goes out of scope
```

Threads (2)

```
There are two solutions:
void fun()
  thread BackgroundPrinting(PrintTextFile, string("c:\\temp\\data.txt"));
  BackgroundPrinting.join(); // main thread waits until the PrintTextFile returns
  return;
void fun()
  thread BackgroundPrinting(PrintTextFile, string("c:\\temp\\data.txt"));
  BackgroundPrinting.detach(); // we get a thread that runs independently.
 return;
```

If a thread is detached (called also as daemon thread), it just keeps running in the background until the end of entry point function. Destroying of the thread object does not stop the detached thread. However, when the process dies (i.e. the *main()* returns), it kills all the running detached threads.

If a thread is not detached and we destroy the thread object without applying *join*, we'll get a run-time error.

Threads (3)

If the entry point function is a member of a class: thread thread object name(&class name::entry point function name, pointer to object, list of input parameters); Example: suppose we have class class Printers PrintTextFile(string FileName); and object: Printers *pHP LaserJet = new Printers; then to create and launch the background printing write: thread BackgroundPrinting(&Printers::PrintTextFile, pHP laserJet, string("c:\\temp\\data.txt"));

If the function launching thread and the entry point function are members of the same class and belong to the same object, *pointer_to_object* is of course pointer *this*.

Threads (4)

```
The entry point may be specified by lambda expressions:
thread thread object name(lambda expression);
Example:
const char *pTextFile = "c:\\temp\\data.txt";
thread BackgroundPrinting([&]() { PrintTextFile(pTextFile); });
Also, it is possible to write a new class and apply functors. Example:
class PrintThread
private: const char *pFile;
public: PrintThread(const char *p) : pFile(p) { } // parameter is specified in constructor
        operator() () const { PrintTextFile(pFile); }
PrintThread pt("c:\\temp\\data.txt"); // create functor object
thread BackgroundPrinting(pt); // functor object presents the entry point function
or
thread BackgroundPrinting(PrintThread("c:\\temp\\data.txt")); // create nameless functor
or with uniform initialization
thread BackgroundPrinting { PrintThread("c:\\temp\\data.txt") }; // here we create object
// BackgroundPrinting and initialize it with nameless functor
```

Threads (5)

```
Alternative:
class PrintThread
{
    public: // no constructor, just operator()
        operator() (const char *p) const { PrintTextFile(p); }
};
thread BackgroundPrinting(PrintThread(), "c:\\temp\\data.txt"); // create nameless object
or
PrintThread pt;
thread BackgroundPrinting(pt, "c:\\temp\\data.txt"); // use previously defined object
```

Threads (6)

If the thread object is destroyed when the thread is still running, the program will creash. The solution is to use method *join*:

Worker.join(); // blocks function Test() until the end of thread

The problem here is that function *Test* may have several *return* amd *throw* points, so we need to call *join* in many places. If, in addition, we have more than one thread, we may get very complicated and indistinct code.

If we instead of class *thread* apply C++ v. 20 class *jthread*, method *join* is called automatically by the destructor of this class and the code presented on this slide will work.

Threads (7)

```
void Test () {
      .....// Do something
jthread Worker(/* entry point function and input parameters */);
 ......// Do something
 if (/* some condition */) {
   return; // no need to call join
 try {
     ......// Do something
  catch (exception) {
   return; // no need to call join
Worker.join();
  ......// the thread has finished, analyse the results
```

jthread provides the same interface as *thread* so it is possible in old code simply add character j to classname *thread*. Header file *<thread>* presents the both classes.

Class this_thread

Class *this thread* has several static methods allowing the current thread to control itself. this thread::sleep for(duration); blocks the current thread for the specified time interval. Example: this thread::sleep for(chrono::seconds(5)); this thread::sleep until(time point); blocks the current thread until the specified moment. Example: time t now t = chrono::system clock::to time t(chrono::system clock::now()); struct tm now tm; localtime s(&now tm, &now t); now tm.tm sec = now tm.tm min = 0; now tm.tm hour++; this thread::sleep until(chrono::system clock::from time t(mktime(&now tm)));

// for example, if time is now 14:05:00 then the thread resumes its work at 15:00:00

Exceptions in threads (1)

Suppose that our main thread has launched another thread. An exception has occurred in this thread but for some reasons the handling of it in the launched thread itself is not possible and should be transferred to the main thread. For solution the problem we need:

- An object of class *exception_ptr*, used for storing the exception. This object is actually a smart pointer. It must be accessible for the launched thread as well as for the main thread.
- Function *current_exception()* creating from ordinary exception an object of class *exception_ptr* (actually storing the exception).
- Function *rethrow_exception()* allowing to transform the object of class *exception_ptr* back to ordinary exception.

See also https://en.cppreference.com/w/cpp/error/exception ptr.html

Exceptions in threads (2)

```
int main()
 exception ptr Ex ptr = nullptr; // smart pointer Ex ptr will store the exception
 thread Task(TaskMain, &Ex ptr, .....);
 Task.join(); // wait for end of thread Task
 try
    if (Ex ptr)
    { // was there an exception in the thread?
        rethrow exception(Ex ptr); // transform back to normal exception
 catch (const exception& e)
   // process the restored exception
     cout << e.what() << endl;</pre>
 return 0;
```

Exceptions in threads (3)

In thread:

```
void TaskMain(exception_ptr *pEx_ptr, .....)
 try
 catch (exception)
   *pEx ptr = current exception(); // instead of handling store the exception
   return;
```

Race conditions (1)

Mostly, threads share some resources. Suppose we have two threads and they share an integer *static int* x = 1; One of them increments and the other decrements it. Those operations are performed in 3 steps: retrieving the value from memory, incrementing or decrementing and sending back into the memory. As the thread scheduling mechanism can swap between threads at any time, the final result is unpredictable:

Thread 1	Thread 2
Retrieve x (value 1)	
Increment x	
Store x (value 2)	
	Retrieve x (value 2)
	Decrement x
	Store x (value 1)

Thread 1	Thread 2
Retrieve x (value 1)	
	Retrieve x (value 1)
	Decrement
	Store x (value 0)
Increment	
Store x (value 2)	

Thread 1	Thread 2
	Retrieve x (value 1)
	Decrement
	Store x (value 0)
Retrieve x (value 0)	
Increment	
Store x (value 1)	

Thread 1	Thead 2
Retrieve x (value 1)	
Increment	
	Retrieve x (value 1)
	Decrement
Store (value 2)	
	Store (value 0)

Race conditions (2)

From https://support.microsoft.com/en-us/help/317723/description-of-race-conditions-and-deadlocks:

A race condition occurs when two threads access a shared variable at the same time. The first thread reads the variable, and the second thread reads the same value from the variable. Then the first thread and second thread perform their operations on the value, and they race to see which thread can write the value last to the shared variable. The value of the thread that writes its value last is preserved, because the thread is writing over the value that the previous thread wrote.

For us the thread execution is nondeterministic, therefore we cannot control the time or order of execution. It may happen that the "race" does not occur at all (tables 1 and 2). But if it does occur, the "winner" may be any of the threads and therefore the final result is unpredictable (tables 3 and 4).

The solution is to lock the shared memory. It means that when one of the threads is operating with the shared memory field, the other thread(s) must stop and wait for the first one to release the shared memory.

Race conditions (3)

More generally, we have to synchronize the threads:

- When one of the threads owns a resource (array, linked list, disk file, COM port, etc.), the other threads can access this resource for reading only. The have no right to change anything on that resource.
- If a thread needs a resource owned by another thread for writing (i.e. changing data), it must wait until the resource is released.
- When the thread owning a resource does not need it more, it has to release the resource and signal about it to the other threads.

Key problems: communication between threads, locking resources.

The code must be thread-safe.

Mutexes (1)

Let us have two threads:

- 1. The producer thread retrieves data from an external source or generates data in some other way. If a package of data is ready, it stores it in a buffer. Time interval needed to create a package is occasional.
- 2. The consumer thread reads the data from the buffer it shares with the producer and processes it and / or views on the screen. Time interval needed to process a package and view it is also occasional.

It may happen that the consumer starts to process data when the producer has not yet finished the storing of new package. It may also happen that the producer starts to store the new package when the consumer has not yet finished the processing of previous package.

Consequently, to synchronize the work of threads we need a mechanism guaranteeing that

- 1. When the producer is accessing the shared buffer, the consumer has no right to access it. If the consumer tries to access the buffer, it will be blocked until the producer releases the buffer.
- 2. When the consumer is accessing the shared buffer, the producer has no right to access it. If the producer tries to access the buffer, it will be blocked until the consumer releases the buffer.

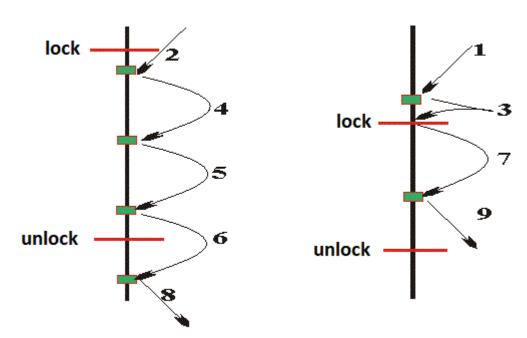
In C++ this mechanism is implemented by mutexes (mutual exclusion classes).

Mutexes (2)

#include <mutex> // see https://en.cppreference.com/w/cpp/thread/mutex.html mutex mx; // must be accessible in both threads

```
// Producer:
mx.lock();
// code in which buffer is accessed
mx.unlock();
```

// Consumer:
mx.lock();
// code in which buffer is accessed
mx.unlock();



Actually, we do not lock the shared memory field but mark the section of code in which the shared memory is accessed.

If the threads need the shared memory only for reading, mutexes are unnecessary.

Mutexes (3)

```
class Producer
{ // task: create random data on an occasional moment
public:
 vector<int> *pBuf; // shared resource
 mutex *pMx;
 Producer(vector<int> *pv, mutex *pm) : pBuf(pv), pMx(pm) { } // constructor
 void operator() ()
distr
    uniform int distribution<int> value distribution(0, 100);
    pMx->lock();
    this thread::sleep for(chrono::milliseconds(delay distribution(generator)));
          // pause, random duration
    ranges::generate(*pBuf, [&]() { return value distribution(generator);});
          // fill the vector with random numbers
    pMx->unlock();
```

Mutexes (4)

```
class Consumer
{ // task: print the created data
public:
 vector<int> *pBuf; // shared resource
 mutex *pMx;
 Consumer(vector<int> *pv, mutex *pm) : pBuf(pv), pMx(pm) { } // constructor
 void operator() ()
   this thread::sleep for(chrono::milliseconds(2500)); // see comment on the next slide
   pMx->lock();
   for each(pBuf->begin(), pBuf->end(), [](int i) { cout << i << ' '; });
   cout << endl;
   pMx->unlock();
};
```

Mutexes (5)

Comment:

The producer starts first but we cannot be sure that it locks the mutex before the consumer applies the locking. If the consumer is faster, it prints zeroes and then allows the producer to work. Therefore we have blocked the consumer for a while. This is not a good way to solve the problem. Later we'll see how to solve it with conditional variables.

Mutexes (6)

Method *try lock* works in the following way:

- If the mutex is not locked by another thread, it will be locked
- If the mutex is locked by another thread, it returns immediately with *false*.

```
Usage example:
```

```
while (true)
  if (mx.try lock())
    ......// critical section operations
    mx.unlock();
    break;
  else
```

Mutexes (7)

If a thread is in critical section and hangs, the mutex is never unlocked and the other threads cannot continue. The solution is *timed_mutex* that has methods *try_lock_for* and *try_lock_until*:

- If the mutex is not locked by another thread, it will be locked immediately.
- If the mutex is locked by another thread, it tries to lock it until the specified time interval has elapsed or the specified time point reached. If the locking is still impossible, *false* is returned.

Usage example:

Timed mutexes can operate as ordinary mutexes, i.e. they support also methods *lock()* and *try lock()*.

Mutexes (8)

To simplify locking and unlocking with mutexes use class lock_guard lock_guard_name(associated_mutex_name);

Example:

#include <mutex> // see https://en.cppreference.com/w/cpp/thread/lock_guard.html mutex mx;

lock_guard<mutex> lock(mx); // constructor for object lock of class lock_guard

The *lock guard* has only two methods: contructor and destructor.

- When the constructor is called, the *lock()* method of the associated mutex is also called.
- When the destructor is applied, the *unlock()* method of the associated mutex is also called.

With manual locking, you have to ensure that the mutex is unlocked correctly on every exit path from the region where you need the mutex locked, including when the region is exited due to an exception. Having a local *lock_guard* object (i.e. auto memory class) it is not a problem: even if an exception is thrown, the destructor is applied and the mutex released. The mutex itself may be a global variable or input parameter of the thread entry point function.

Mutexes (9)

The unique_lock is more flexible:
unique_lock<mutex> unique_lock_name(associated_mutex_name);

// as lock_guard – the constructor locks the associated mutex.
unique_lock<mutex> unique_lock _name(associated_mutex_name, std::defer_lock);

// the constructor does not lock the associated mutex
To lock later you may use methods:

- *lock()* as the *lock()* of *mutex*
- *try lock()* as the *try lock()* of *mutex*
- *try_lock_for* as the *try_lock_for* of *timed_mutex*
- try lock until as the try lock until of timed mutex

The destructor of *unique_lock* unlocks the associated mutex. But you may also unlock with *unlock()* method. Example:

```
mutex mx; // see https://en.cppreference.com/w/cpp/thread/unique_lock.html
unique_lock<mutex> lock (mx, std::defer_lock);
.....
lock.lock(); // locks the mutex mx
.....
if (...)
    return; // unlocks mutex mx
.....
lock.unlock(); // unlocks the mutex mx
```

Mutexes (10)

Call a function just once (1)

Suppose we have two threads sharing a common but not initialized yet resource. The code of both threads starts with the call to initializing function. For example, one of the threads sends data to an external device, the other receives data and the shared device is a COM port or a TCP/IP socket. It is unknown which of the threads will call the initializer first. But it is clear that the initializer must be called only once.

```
void ::initializer();
```

```
// initializes the common resource, for simplicity let it to be out of classes class Thread_1 { .... void operator()() { initializer(); ... } .....}; class Thread_2 { .... void operator()() { initializer(); ... } .....};
```

The most effective solution to similar problems is to apply the *call_once* mechanism:

- 1. Define variable of type *once_flag*, it must be accessible for both threads: once flag flag name;
- 2. In both threads call the initializer in this way: call once(flag, pointer to initializer, list of initializer parameters);

If the first thread wants to call the initializer when the second thread is currently intializing, it blocks until the end of initializing and then continues without stepping into initializer. If the first thread wants to call the initializer when the second thread has already finished the initializing, it continues immediately. Thus the faster thread performs the initialization and the slower omits this step.

Call a function just once (2)

If you have several functions to be called only once you need also several flags.

The initializer may be a function, class method or functor. About details see https://en.cppreference.com/w/cpp/thread/call_once.html/.

```
Example:
```

```
#include <mutex>
class Calculate 1
public:
 once flag *pFlag;
 fstream *pSharedFile;
 Calculate 1(once flag *p1, fstream *p2) : pFlag(p1), pSharedFile(p2) { }
 void operator() ()
    call once(*pFlag, OpenFile, "c:\\temp\\data.txt", &pSharedFile);
// class Calculate 2 is similar
```

Call a function just once (3)

```
int main()
{
  once_flag flag;
  fstream *pDataFile;
  thread thread_1{ Calculate_1(&flag, pDataFile) };
  thread thread_2{ Calculate_2(&flag, pDataFile) };
  thread_1.join();
  thread_2.join();
  return 0;
}
```

One of the threads (impossible to say which) opens the file

Atomic variables (1)

If a single variable is shared, we may lock the operations with it using mutexes. But it is more efficient to declare them as atomic variables.

An atomic operation is a machine instruction that is always executed without interruption. A sequence of two or more machine instructions isn't atomic since the operating system may suspend the execution of the current sequence of operations in favour of another task. A C++ statement is atomic if the compiler translates it into a single machine instruction. But each hardware architecture may translate the same C++ statement in its own different way. So it is wise to say that none of the C++ is statements is atomic. See also slide *Race conditions* (1).

An operation with atomic variables of course needs a sequence of machine instructions. But the execution of this sequence is not interrupted. In other words, an operation with atomic variables is locked although we need not to declare any mutexes for them.

```
#include <atomic> // see <a href="https://en.cppreference.com/w/cpp/atomic/atomic.html">https://en.cppreference.com/w/cpp/atomic/atomic.html</a>
using namespace std;
atomic<int> i(0);
atomic<char> c('A');
atomic<long> l(100000);
........................// all the integral types are allowed
i, c, l etc. are objects. Their initialization is compulsory.
```

Atomic variables (2)

Due to operator functions we may write simply:

```
atomic<int> i = 0;
atomic<char> c = 'A';
atomic<long> l = 100000;
```

Functions declared in atomic classes are interlocked functions. It means that if one of the functions associated with an atomic variable is running and another thread tries to call this or some other function of the same atomic variable, this thread has to wait. In other words, the functions of atomic classes lock and unlock their inner contents automatically. For example:

```
i.store(1);  // or due to operator functions we may write also i = 1; if (i.load() == 1) .....  // or due to operator functions we may write also if (i == 1)

Some examples about the other atomic class functions:
i.fetch_add(1);  // or due to operator functions we may write also i++;
i.fetch_add(10);  // or due to operator functions we may write also i+=10;
i.fetch_sub(1);  // or due to operator functions we may write also i--;

But:
atomic<int> i1 = 10;
atomic<int> i2 = i1; // error
i2 = i1; // error
i2.exchange(i1); // meaning: i2 = i1
```

Atomic variables (3)

Similarly:

```
atomic<int> i1 = 10, i2 = 20, i3 = 0;
i3.exchange(i1 + i2); // i3 = i1 + i2
cout << i3 << endl;
i3.exchange(i1 * i2); // i3 = i1 * i2
Atomic pointers are also allowed, for example:
atomic<vector<int> *> pv = new vector<int>(10);
atomic<int *> pi = new int[10];
```

Pointer operations are supported with overloaded operators, for example *pi = 10;

$$*(pi + 1) = 20;$$

Thread interrupt (1)

Rather often a thread runs in endless loop and we need a mechanism to interrupt it. In C++ killing of thread is not allowed. What we need are tools for so called cooperative (sometimes said polite) cancellation with which we force the thread entry point function to return. In the following example we use an *atomic*<*bool*> variable for that:

```
class Worker {
public:
     vector<int> &buf;
     atomic<br/>bool> &stop;
     Worker (vector<int> &v, atomic<bool> &b) : buf(v), stop(b) { }
    void operator() () {
      default random engine generator;
      uniform int distribution<int> delay distribution(0, 10000);
      uniform int distribution<int> value distribution(0, 100);
       while (!stop) { // worker runs until stop becomes true (polls value of variable stop)
         this thread::sleep for(chrono::milliseconds(delay distribution(generator)));
         ranges::generate(buf, [&]() { return value distribution(generator); });
         ......// do something with the generated numbers
```

Thread interrupt (2)

```
class Controller
public:
       atomic<br/>bool> &stop;
       Controller(atomic < bool > &b) : stop(b) { }
       void operator() ()
          getch(); // the human operator has to press a key to stop worker
          stop = true;
int main()
      atomic<br/>bool> stop = false;
      vector<int> buf(32);
      thread ControllerThread { Controller(stop) };
      thread WorkerThread { Worker(buf, stop) };
      WorkerThread.join();
      ControllerThread.join();
      return 0;
```

Thread interrupt (3)

In the previous example the thread entry point function checks periodically the stop flag and quits if an outer function has set it to *false*. Class *jthread* has better tools:

```
#include <stop token> // see <a href="https://en.cppreference.com/w/cpp/header/stop">https://en.cppreference.com/w/cpp/header/stop</a> token
void Worker(int n, stop token stop) { // obligatory parameter of type stop token
  for (int i = 0; i < n \&\& !stop.stop requested(); <math>i++) {
    this thread::sleep for(chrono::seconds(1));
    cout << i << endl;
  } // the thread must time to time poll the state of stop token object
void Test {
  stop source; // to control the interrupting
  stop token stop = source.get token();
  jthread thr(Worker, 10, stop);
  this thread::sleep for(chrono::seconds(5)); // allow to run 5 seconds and then interrupt
  source.request stop(); // after that stop requested() in Worker returns true
```

Thread interrupt (4)

It is possible to associate the stop token object with a callback (function, functor or lambda). The callback is invoked when the stop is requested: void Worker(int n, stop token stop) stop callback OnRequest(stop, []() { cout << "Broken" << endl; });</pre> for (int i = 0; i < n && !stop.stop requested(); <math>i++) {// the thread must time to time check the state of stop token object this thread::sleep for(chrono::seconds(1)); $cout \ll i \ll endl;$ void Test stop source source; // to control the interrupting stop token stop = source.get_token(); jthread thr(Worker, 10, stop); this thread::sleep for(chrono::seconds(5)); source.request stop(); // after that stop requested() in Worker returns true // in addition callback OnRequest is invoked, // i.e. message "Broken" is printed

Condition variables (1)

Inter-thread communication in C++ is implemented by condition variables. Let us have two threads:

- 1. The producer thread retrieves data from an external source or generates data in some other way. If a package of data is ready, it stores it in a buffer. Time interval needed to create a package is occasional.
- 2. The consumer thread reads the data from the buffer shared with the producer and processes it and / or views on the screen. Time interval needed to process a package and view it is also occasional. Problem: the processing may start only when the producer has finished its task.

Condition variables (2)

wait(unique lock name);

The *wait* method can be called only when the thread is already locked by *unique_lock* managing the associated mutex. It:

- 1. Automatically releases the lock (atomically calls *unlock()* of the associated *unique_lock*), thus allowing the other thread to run.
- 2. Blocks its own thread and starts to wait for a notification.
- 3. Unblocks its thread when the other thread has emitted the notification.
- 4. Locks again (before return atomically calls *lock()* of the associated *unique lock)*.

Notification is emitted by method *notify_one()* called from another thread.

If several threads are waiting, *notify_one()* releases only one of them (impossible to say which). Method *notify_all()* sends notification to all the threads stopped by the current conditional variable. If there are no waiting threads, those functions do nothing.

Condition variables (3)

wait(unique_lock_name, predicate);

The predicate may be a pointer to function, lambda expression or functor. It cannot have any arguments and it must return *true* or *false*.

wait() with predicate:

- 1. If the predicate has returned *true*, does nothing and returns immediately.
- 2. If the predicate has returned *false*, releases the lock (atomically calls *unlock()* of the associated *unique_lock*), thus allowing the other thread to run.
- 3. Also, if the predicate has returned *false* blocks its thread and starts to wait for a notification.
- 4. If the notification has arrived, calls the predicate once more.
- 5. If the predicate still returns *false*, the thread stays in blocked state.
- 6. If the predicate returns *true*, unblocks its thread and locks the associated *unique_lock*.

std::notify_all_at_thread_exit(condition_variable_name, unique_lock_name);

This function not belonging to any classes is mostly called before *join* to release threads that may be have got stuck.

Condition variables (4)

```
Example:
int main()
 vector<int> buf; // empty vector
 mutex mx;
 condition variable cv;
 thread ProducerThread { Producer (buf, mx, cv, 32) };
 thread ConsumerThread { Consumer (buf, mx, cv) };
 ProducerThread.join();
 ConsumerThread.join();
 return 0;
// the order of launching threads is meaningless
```

Condition variables (5)

```
class Producer
public:
  vector<int> &buf;
  mutex &mx;
  condition variable &cv;
  int buf size;
  Producer (vector<int> &v, mutex &m, condition variable &c, int n):
            buf(v), mx(m), cv(c), buf size(n) { }
 void operator() ()
     default random engine generator;
     uniform int distribution<int> delay (0, 10000);
     uniform int distribution<int> random number(0, 100);
     unique lock<mutex> lock(mx); // locks the mutex
     buf.resize(buf size);
     this thread::sleep for(chrono::milliseconds(delay(generator)));
     ranges::generate(buf, [&]() { return random number(generator); });
     cv.notify one(); // releases the comsumer
  } // variable "lock" deleted, mutex automatically unlocked
```

Condition variables (6)

```
class Consumer
public:
 vector<int> &buf;
 mutex &mx;
 condition variable &cv;
  Consumer(vector<int> &v, mutex &m, condition variable &c):
             buf(v), mx(m), cv(c) { }
 void operator() ()
      unique lock<mutex> lock(mx); // locks the mutex
      cv.wait(lock, [&]() { return !buf.empty(); });
            // If the buffer is empty, unlocks the mutex allowing the producer to work
            // and starts to wait for the notification from the producer.
            // The producer sends notification when the buffer is filled. Consequently
            // predicate then returns true and the thread will be unblocked.
      ranges::for each(buf, \lceil \rceil (int i) { cout \leq i \leq ' '; });
      cout << endl;
}; // destructor called, mutex automatically unlocked
```

Condition variables (7)

Comments:

- 1. Suppose that the producer steps into the critical section first. In that case the consumer must wait until the producer has left the critical section. i.e. until the end of producer thread. The first step of the consumer is the call to *wait()*. As the buffer is already full, the predicate returns *true* and *wait()* returns immediately the consumer may start to run. The notification sent by the producer is ignored.
- 2. Now suppose that the consumer steps into the critical section first and thus blocks the producer. The first step of the consumer is the call to *wait()*. As the buffer is empty, the predicate returns *false*. Now:
 - a. wait() unlocks the mutex, the producer may start to run.
 - b. wait() does not return but starts to wait for the notification. So the consumer is blocked and the producer may complete its task.

Before return the producer sends notification. *wait()* in consumer calls once more the predicate. As the buffer is already filled, the predicate returns *true*. Now:

- a. wait() locks the mutex, but as the producer has already finished, it has no consequences.
- b. wait() returns and thus the consumer may start to run.

So, in this solution we may be sure that at first the producer creates the data and only after that the consumer processes the data. Compare it with solution on slides *Mutexes* (3) ... *Mutexes* (5).

Condition variables (8)

Example: classical producer – consumer problem in which the producer has to inform the the consumer that a data package is prepared, wait until the consumer has finished the processing and start to prepare the next package.

```
int main()
 atomic<br/>bool> stop = false;
 vector<int> *pBuf = new vector<int>;
 mutex mx;
 condition variable cv;
 thread ControllerThread { Controller(stop) }; // see slide Atomic variables (5)
 thread ProducerThread { Producer(pBuf, &stop, &mx, &cv, 32) };
 thread ConsumerThread { Consumer(pBuf, &stop, &mx, &cv) };
 ConsumerThread.join();
 ProducerThread.join();
 ControllerThread.join();
 return 0;
```

Condition variables (9)

```
class Producer {
public: vector<int> *pBuf;
       atomic<bool> *pStop;
       mutex *pMx;
       condition variable *pCv;
       int buf size;
       Producer (vector<int> *pv, atomic<bool> *pb, mutex *pm, condition variable *pc,
                 int n): pBuf(pv), pStop(pb), pMx(pm), pCv(pc), buf size(n) { }
       void operator() () {
            default random engine generator;
            uniform int distribution<int> delay (0, 10000);
            uniform int distribution<int> random number(0, 100);
            while (!*pStop) {
                unique lock<mutex> lock(*pMx); // locks the mutex
                pBuf->resize(buf size);
                this thread::sleep for(chrono::milliseconds(delay(generator)));
                ranges::generate(*pBuf, [&]() { return random number(generator); });
                pCv->notify one(); // releases the comsumer
                pCv->wait(lock); // blocks the producer
```

Condition variables (10)

```
class Consumer {
public: vector<int> *pBuf;
       atomic<bool> *pStop;
       mutex *pMx;
       condition variable *pCv;
       Consumer (vector<int> *pv, atomic<bool> *pb, mutex *pm, condition variable *pc) :
                   pBuf(pv), pStop(pb), pMx(pm), pCv(pc) { }
       void operator() () {
            while (!*pStop) {
                unique lock<mutex> lock(*pMx); // locks the mutex
                pCv->wait(lock, [&]() { return !pBuf->empty(); });
                ranges::for each(*pBuf, \lceil \mid (int i) \mid (cout << i << ' '; \});
                pBuf->resize(0);
                cout << endl;
                pCv->notify one(); // releases the producer
```

Condition variables (11)

```
wait_for(unique_lock_name, duration);
wait_until(unique_lock, time_point);
and
wait_for(unique_lock_name, duration, predicate);
wait_until(unique_lock_name, timepoint, predicate);
```

Those two methods block the thread until notified or until the specified timeout has elapsed or timepoint has been reached. The return value is $cv_status::timeout$ or $cv_status::no_timeout$. About the usage of predicates see slide $Conditional\ variables\ (3)$.

Condition variables (12)

It may happen that when the stop is requested (see slide *Thread interrupt (3)*), a thread is waiting for notification for a *condition_variable* and cannot check the state of *stop_token*. The following example presents the solution of this problem:

```
void Test {
 mutex mx;
 condition variable any cv; // use instead of condition variable
                    // see https://en.cppreference.com/w/cpp/thread/condition variable any
 queue<int> values;
 stop source source; // to control the interrupting
 stop token stop = source.get token();
 jthread thr1(Consumer, &mx, &cv, &values, stop); // on the next slides
 jthread thr2(Producer, &mx, &cv, &values, stop);
 this thread::sleep for(chrono::seconds(5)); // allow to run 5 seconds and then interrupt
 source.request stop();
 thr1.join();
 thr2.join();
```

condition_variable can wait only on unique:lock<mutex>. condition_variable_any can wait on any lockable object (i.e. on anything that has lock() and unlock() methods). In addition, with condition_variable_any we can stop waiting not only with notification but also with interrupting the thread.

Condition variables (13)

```
void Producer(mutex* pmx, condition variable any* pcv, queue<int>* pvalues,
              stop token stop)
 // the Producer must insert into a queue 10 random numbers
  default random engine generator;
  uniform int distribution<int> random number(0, 100);
  unique lock<mutex> lock(*pmx);
  for (int i = 0; i < 10; i++) {
    if (stop.stop requested()) {
       cout << "Stop request detected" << endl;</pre>
       while (!pvalues->empty()) {
         pvalues->pop(); // interrupt, clean the queue and exit
       return:
     pvalues->push(random number(generator));
     this thread::sleep for(chrono::milliseconds(1000));
  pcv->notify one(); // queue is filled, allow the consuer to start
  cout << endl;
```

Condition variables (14)

```
void Consumer(mutex *pmx, condition_variable_any *pcv, queue<int> *pvalues,
               stop token stop) {
  unique lock<mutex> lock(*pmx);
  pcv->wait(lock, stop, [&]() { return !pvalues->empty(); });
  // waiting ends when the notification has arrived or when the stop is requested
  if (stop.stop requested()) {
     cout << "Waiting broken off" << endl;</pre>
    return;
  while (!pvalues->empty()) { // prints the results
     cout << pvalues->front() << ' ';</pre>
    pvalues->pop();
  cout << endl;
```

async and futures (1)

Suppose our program consists of 3 tasks. The final task 3 needs data prepared by tasks 1 and 2. There are no any dependences between tasks 1 and 2: they do not need request data from each other, do not share resources, etc. If we have a multiprocessor computer, we may put task 2 into a separate background thread and force it run concurrently with task 1: thus we may get better performance.

```
There is another way: use function std::async and futures:
future <entry point function return_value_type> future_name =
                              async(entry point function name, list of input parameters);
The task entry point function may have any set of input parameters. Example: suppose the
entry point function for task 2 is:
bool ReadData(unsigned char *, int);
To start asynchronous reading task:
#include <future> // see <a href="https://en.cppreference.com/w/cpp/thread/future.html">https://en.cppreference.com/w/cpp/thread/future.html</a>
                   // see <a href="https://en.cppreference.com/w/cpp/thread/async.html">https://en.cppreference.com/w/cpp/thread/async.html</a>
unsigned char *pBuf = new unsigned char[1024 * 10];
future < bool > result = async(ReadData, pBuf, 1024 * 10); // task 2 is now in a separate thread
......// program continues with task 1
if (result.get()) // method get() returns the return value of task entry point function
{ // you can call get on a specific future only once!
 ......// program starts to execute task 3
```

async and futures (2)

async does not guarantee that a separate thread with the specified entry point function will start immediately. When the future's *get()* method is called, it is possible that:

- The asynchronous task has finished and *get()* returns the output value immediately.
- The asynchronous task is still running. In that case *get()* blocks the current thread until the return value has become available.
- The asynchronous task has not started yet. Then it is forced to start, the current thread blocks. Actually it means that the ansynchronous executing of tasks has failed. If the call to *get()* is omitted, it may happen that the asynchronous task never starts.

There are some possibilities to control the behavior of asynchronous task:

```
future< entry_point_function_return_value_type> future_name =
    async(launch_policy, entry_point_function_name, list_of_input_parameters);
```

The launch policy may be:

- 1. launch::async try to launch a new thread immediately. May be risky because if it is not possible, exception will be thrown.
- 2. *launch::deferred* launch when *get()* is called, i.e. no concurrent executing. May be useful as temporary setting for debugging.
- 3. launch::async | launch::deferred default setting, launch time is set by system.

async and futures (3)

Class future has method:

```
wait();
```

This function behaves similarly to get(): if the asynchronous task has not finished it blocks the current thread until the asynchronous task has ended. Also, if necessary it forces the asynchronous task to start. Later, the output value may be retrieved by get().

Class *future* has also methods:

- wait for(duration);
- wait_until(timepoint);

Those functions block the current thread until the asynchronous task has finished or timeout elapsed or the specified timepoint reached. But they do not force the asynchronous task to start. Their return value may be:

- 1. future_status::deferred the asynchronous task has not started yet.
- 2. future_status::timeout the asynchronous task is running but waiting time has elapsed.
- 3. future status::ready the asynchronous task has finished.

Example:

```
if (result.wait_for(chrono::seconds(60)) != future_status::ready) {
    cout << "It seems that the asynchronous task has problems" << endl;
    return;
}</pre>
```

wait_for(chrono::seconds(0)) returns us the current status of the asynchronous task.

async and futures (4)

If the asynchronous task has thrown an unhandled exception, get() catches it and rethrows:

Method wait() does not rethrow exceptions.

Method *get()* may be called only once. After get() the future becomes invalid. To check the state of future use method *valid()* – it returns false if the future has become not useable.

In a complicated application with a lot of threads multiple calls to *get()* may be needed. In that case use *shared_future* (https://en.cppreference.com/w/cpp/thread/shared_future.html).

async and futures (5)

Instead of entry point function we may use functor or lambda. Example:

```
class Producer
public:
  list<int> &buf; // reference
  int lower, upper;
  Producer(list<int> &v, int l, int u) : buf(v), lower(l), upper (u) { }
  void operator() ()
       default random engine generator;
       uniform int distribution<int> delay(0, 10000);
       uniform int distribution<int> random number(lower, upper);
        this thread::sleep for(chrono::milliseconds(delay(generator)));
        ranges::generate(buf, [&]() { return random number(generator); });
```

async and futures (6)

```
int main()
   list<int> buf1(10);
   list<int> buf2(10);
   future < void > f1 = async(Producer(buf1, 0, 100));
   future < void > f2 = async(Producer(buf2, -100, 0));
   fl.get(); // here the thread returns nothing but we need to wait until it quits
             // so get() from class future behaves as join from class thread
   f2.get();
   buf1.sort();
   buf2.sort();
   buf1.merge(buf2);
   ranges::for each(buf1, [](int i) { cout << i << ' '; });
   cout << endl;
   return 0;
```

Shortly: a *future* and *async* provide facilities to retrieve values and / or exceptions from a function that is located in background thread and is currently executing or has already executed or will be executed (in the last case, they may force to start executing).

async and futures (7)

```
Example (compare with solution on slides Thread interrupt (1) ... (4)):
class Worker
public:
vector<int> &buf;
future<void> &fut;
Worker (vector<int> &v, future<void> &f): buf(v), fut(f) { }
void operator() ()
 default random engine generator;
 uniform int distribution<int> delay distribution(0, 10000);
 uniform int distribution<int> value distribution(0, 100);
 while (fut.wait for(chrono::seconds(0)) != future status::ready) {
    // worker runs until the keyboard async thread has not exited
    this thread::sleep for(chrono::milliseconds(delay distribution(generator)));
    ranges::generate(buf, [&]() { return value distribution(generator); });
```

async and futures (8)

```
void Keystroke()
   _getch(); // waits for keystroke, then exits
int main()
   vector<int> buf(32);
   future<void> fut = async(Keystroke);
   thread WorkerThread { Worker(buf, fut) };
   WorkerThread.join();
   return 0;
```

Packaged tasks

It is possible to prepare tasks beforehand and invoke them later (or not invoke at all if they come out to be unnecessary):

```
packaged task<entry point value return type(list of parameter types)>
task name(entry point function name);
Each packaged task contains a future. To retrieve it:
future<entry point value return type> future name = task name.get future();
To invoke task:
task name(actual parameters list);
and after that apply the associated future's get() to retrieve the result.
Example: suppose the entry point function (functors and lambdas also allowed) for a task is:
bool ReadData(unsigned char *, int);
To create the associated task package:
#include <future> // see <a href="https://en.cppreference.com/w/cpp/thread/packaged">https://en.cppreference.com/w/cpp/thread/packaged</a> task.html
packaged task<bool(unsigned char *, int)> read task(ReadData);
.....// program continues
if (data needed) {
 unsigned char *pBuf = new unsigned char[1024 * 10];
 future<br/>
sool> read result = read task.get future(); // retrieves the future
 read task(pBuf, 10240) // program starts to execute task, if possible then in separate thread
 bool success = read result.get(); // blocks until the end of task
```

Promises (1)

We can store the data created by a thread into a promise and later use future to retrieve it into another thread:

```
promise<task return value type> promise name;
promise name.set value(data to store); // can be called only once
Each promise has a future. When we need data stored in promise, retreive it:
future<return_value_type> future_name = promise name.get future();
and then call get().
Example:
void Producer(list<int> *pBuf, int lower, int upper, promise<list<int> *> *pPromise)
{ // see <a href="https://en.cppreference.com/w/cpp/thread/promise.html">https://en.cppreference.com/w/cpp/thread/promise.html</a>
 default random engine generator;
 uniform int distribution<int> delay(0, 10000);
 uniform int distribution<int> random number(lower, upper);
 this thread::sleep for(chrono::milliseconds(delay(generator)));
 ranges::generate(*pBuf, [&]() { return random number(generator); });
 pBuf->sort();
 pPromise ->set value(pBuf); // store the result into promise
                                // turn attention that the thread entry point function
                                 // does not return any value
```

You can call set_value on a specific promise only once!

Promises (2)

```
void Consumer(promise<list<int>*>* pPromise)
  future<list<int> *> fut = pPromise->get future(); // future associated with promise
  list<int> *pRes = fut.get(); // retrieve the result, if necessary, wait
  ranges::for each(*pRes, [](int i) { cout \leq i \leq ' '; });
  cout << endl;
int main()
{ // Compare with similar example on slides Mutexes (3) ... Mutexes (5)
  // Here we need neither mutexes nor shared memory fields
 promise<list<int> *> prom;
 list<int> buf(10);
 thread thr1(Producer, &buf, 0, 100, &prom);
 thread thr2(Consumer, &prom);
 thr1.join();
 thr2.join();
 return 0;
```

Promises (3)

The promises are very necessary if we need to get results created by detached threads. Example:

```
void Producer(list<int> *pBuf, int lower, int upper, promise<void> *pPromise)
 default_random engine generator;
 uniform int distribution<int> delay(0, 10000);
 uniform int distribution<int> random number(lower, upper);
 this thread::sleep for(chrono::milliseconds(delay(generator)));
 ranges::generate(*pBuf, [&]() { return random number(generator); });
 pBuf->sort();
 pPromise ->set value(); // inform that the thread has ended
void Consumer(list<int> *pBuf, promise<void>* pPromise)
  future<void> fut = pPromise->get future(); // future associated with promise
  fut.get(); // wait until the end of producer
  ranges::for each(*pRes, [](int i) { cout << i << ' '; });
  cout << endl;
```

Promises (4)

```
int main()
{
  promise<void> prom;
  list<int> buf(10);
  thread thr1(Producer, &buf, 0, 100, &prom);
  thr1.detach();
  thread thr2(Consumer, &buf, &prom);
  thr2.join();
  return 0;
}
```

Latches (1)

```
C++ v. 20 has some new tools for synchronization of threads. One of them is the latch:
#include <latch> // see <a href="https://en.cppreference.com/w/cpp/thread/latch">https://en.cppreference.com/w/cpp/thread/latch</a>
void Test () {
  vector<int>v1, v2;
  latch data ready(2), clear data(1);
  // Latch has a counter, its initial value is the parameter of constructor
  // There is no possibility to increase or reset the value of counter later
  jthread thr1(Producer1, &v1, 10, 500, &data ready, &clear_data); // see the next slide
  jthread thr2(Producer2, &v2, 5, 200, &data ready, &clear data);
  data ready.wait(); // Wait until the latch counter is zero. The threads decrement the
                       // counter. So, Test() can continue when the both vectors are filled
  ranges::sort(v1);
  ranges::sort(v2);
  vector\leqint\geq v(15);
  ranges::merge(v1, v2, v.begin());
  ranges::for_each(v, [](const int& i) { cout << i << ' '; });
  cout << endl;
  clear data.count down(); // decrements the counter, allow the threads to continue
  thr1.join();
  thr2.join();
```

Latches (2)

```
void Producer1(vector<int> *pvec, int n, int t, latch *pdata ready, latch *pclear data) {
 default random engine generator;
 uniform int distribution<int> random number(0, 100);
 for (int i = 0; i < n; i++) { // fill the vector
   pvec->push back(random number(generator));
   this thread::sleep for(chrono::milliseconds(t));
 pdata ready->count down(); // atomically decrements the latch counter
 pclear data->wait(); // waits until latch clear data decremented by Test() becomes 0
 pvec->clear(); // data is consumed, we can now delete it
void Producer2(vector<int>* pvec, int n, int t, latch *pdata ready, latch *pclear data) {
 default random engine generator;
 binomial distribution<int> random number(100);
 for (int i = 0; i < n; i++) {
   pvec->push back(random number(generator));
   this thread::sleep for(chrono::milliseconds(t));
 pdata ready->count down();
 pclear data->wait();
 pvec->clear();
```

Barriers (1)

The barrier is more flexible: #include <barrier> // see https://en.cppreference.com/w/cpp/thread/barrier void Test () vector<int>v1, v2; barrier data ready { 3}; // barrier for 3 tasks // here template barrier has default parameter ithread thr1(Producer1, &v1, 10, 500, &data ready); // see the next slide jthread thr2(Producer2, &v2, 5, 200, &data ready); data ready.arrive and wait(); // The first task is performed: the threads are launched // Waits until the other 2 tasks implemented by Producer1 and Producer2 are performed // When all the 3 tasks are marked as done, stops waiting ranges::sort(v1); ranges::sort(v2); vector \leq int \geq v(15); ranges::merge(v1, v2, v.begin()); ranges::for each(v, [](const int& i) { cout \leq i \leq ' '; }); cout << endl; thr1.join(); thr2.join();

Barriers (2)

```
void Producer1(vector<int>* pvec, int n, int t, barrier<> *pdata ready) {
 default random engine generator;
 uniform int distribution<int> random number(0, 100);
 for (int i = 0; i < n; i++)
   pvec->push back(random number(generator));
   this thread::sleep for(chrono::milliseconds(t));
 pdata ready->arrive and wait();
          // marks the task as done (vector is filled), waits until all the tasks are closed
void Producer2(vector<int>* pvec, int n, int t, barrier<> *pdata ready) {
 default random engine generator;
 binomial distribution<int> random number(100);
 for (int i = 0; i < n; i++) {
   pvec->push back(random number(generator));
   this thread::sleep for(chrono::milliseconds(t));
 pdata ready->arrive and wait();
```

Barriers (3)

The barrier may have a callback function. It is invoked when all the tasks are marked as done. It is implemented as a functor:

```
class Msg {
public: void operator() () noexcept { cout << "Ready" << endl; } // noexcept is necessary
void Producer1(vector<int>* pvec, int n, int t, barrier<Msg> *pdata ready) {
  void Producer2(vector<int>* pvec, int n, int t, barrier<Msg> *pdata ready) {
void Test () {
 barrier Msg > data ready { 3 }; // template parameter is the callback typename
   . . . . . . . . . . . . . . . . . . .
```

Remark: about keyword *noexcept* read slide *C*++ *standard exceptions* (3) from *chapter Cpp standard functions*, course IAX0586

Barriers (4)

When all the tasks are done and the *arrive_and_wait()* method returns, the number of tasks specified in the constructor is reset. Thus, a barrier may be used in loops. The following example presents a producer-consumer problem solution implemented with barriers:

```
void Test ()
 barrier <> data ready { 2 };
 vector<int> v;
 default random engine generator;
 for (int i = 0; i < 5; i++)
   jthread* pthr1 = new jthread { Producer, &generator, &v, 10, &data ready };
   jthread* pthr2 = new jthread { Consumer, &v, &data ready };
   pthr1->join();
   pthr2->join();
   delete pthr1;
   delete pthr2;
```

Barriers (5)

```
void Producer(default random engine *pgen, vector<int>* pvec, int n, barrier<>* pready)
 uniform int distribution<int> random number(0, 100);
 pvec->clear();
 for (int i = 0; i < n; i++)
   pvec->push back(random number(*pgen));
   this thread::sleep for(chrono::milliseconds(500));
 pready->arrive and wait();
void Consumer(vector<int>* pvec, barrier<>* pready)
 pready->arrive and wait();
 ranges::for each(*pvec, [&](const int& i)
     cout << i << ' ';
     this thread::sleep for(chrono::milliseconds(500));
    });
 cout << endl;
```

Semaphores (1)

C++ v. 20 has two types of semaphores:

#include <semaphore> // https://en.cppreference.com/w/cpp/thread/counting_semaphore int max_value, initial_value;

counting_semaphore<max_value> sem1(initial_value);

binary_semaphore sem2(initial_value); // max_value is 1, initial_value may be 0 or 1

Method *release()* atomically increments the counter, method *acquire()* decrements it. The counter cannot be negative and cannot be greater than the *max_value*.

If the counter has become zero, *acquire()* blocks the thread. If due to call to *release()* the counter has a positive value, the blocked thread can continue.

An example about usage of semaphores: suppose we have a server that must process requests. For processing a new request we have to launch a new thread. However, the number of threads running concurrently cannot be endless.

To solve the problem we start the program with semaphore in which the <code>initial_value</code> is set to max. Each thread starts with call to <code>acquire()</code>, i.e. with starting the thread we decrement the counter. If the counter becomes zero, max allowed number of threads are already running and the new thread must wait. Each thread ends with call to <code>release()</code>, i.e. with ending the thread we increment the counter. If the counter was 0, it is now 1 and the thread that was blocked may start to run.

The *binary_semaphore* can replace mutexes. See the example on the next slide.

Semaphores (2)

```
void Producer(vector<int>* pvec, int n, binary semaphore* pdone) {
  default random engine generator;
  uniform int distribution<int> random number(0, 100);
  for (int i = 0; i < n; i++) {
    pvec->push back(random number(generator));
    this thread::sleep for(chrono::milliseconds(500));
  pdone->release();
void Consumer(vector<int>* pvec, binary semaphore* pdone) {
  pdone->acquire(); // blocked until the Producer increments the counter
  ranges::for each(*pvec, [\&](const int& i) { cout << i << ' '; });
  cout << endl;
void Test () {
  binary_semaphore done(0); // initially in state 0
  vector<int> v;
  jthread thr1(Producer, &v, 10, &done);
  jthread thr2(Consumer, &v,&done);
  thr1.join();
  thr2.join();
```

Asynchronous I/O in Windows (1)

C++ I/O standard classes (see https://en.cppreference.com/w/cpp/io/basic_fstream.html are excellent for disk file operations. For reading from and writing into external devices connected over COM port or TCP socket, we have to use the Windows mechanisms.

The first step is to create the file:

```
HANDLE handle name = // HANDLE is defined in Windows.h
CreateFileA(file name and path, // as regular C string
desired access, // GENERIC READ for files used for reading only
               // GENERIC WRITE for files used for writing only
               // GENERIC READ | GENERIC WRITE for the both operations
share mode, // outside the course scope, set to 0
security attributes, // outside the course scope, set to NULL
creation disposition, // CREATE ALWAYS and if already exists, at first destroy it
                   // CREATE NEW and if already exists, the operation fails
                   // OPEN EXISTING and if not found, the operation fails
                   // OPEN ALWAYS and if not found creates a new one
                   // TRUNCATE EXISTING (destroy the contents) and if not found,
                   // the operations fails
flags and attributes, // in our course FILE_FLAG_OVERLAPPED (discussed later)
template file handle); // outside the course scope, set to NULL
```

The complete specification of function *CreateFileA* is on website https://docs.microsoft.com/en-us/windows/desktop/api/fileapi/nf-fileapi-createfilea is

Asynchronous I/O in Windows (2)

If the function fails, the return value is *INVALID_FILE_HANDLE*. To know the reason, call function *GetLastError*:

```
unsigned long error_code = GetLastError();
```

The error codes are on https://docs.microsoft.com/en-us/windows/desktop/Debug/system-error-codes.

If you do need the file any more, close it:

CloseHandle(handle_name);

CloseHandle(hFile);

Example:

```
HANDLE hFile = CreateFileA("FileExample.bin", GENERIC_READ | GENERIC_WRITE, 0, NULL, CREATE_ALWAYS, FILE_FLAG_OVERLAPPED, NULL); if (hFile == INVALID_HANDLE_VALUE) cout << "File not created, error " << GetLastError() << endl; // Important: you should always check was a file operation successful or not ......
```

Asynchronous I/O in Windows (3)

```
To read from file synchronously:
int result name = ReadFile(handle name, pointer to buffer that receives read data,
max number of bytes to read, pointer to variable for number of bytes actually read,
NULL);
If the reading failed, the return value is FALSE and GetLastError() returns the error code.
The actual number of read bytes may be less than the number of needed bytes.
Example:
unsigned long nBytesToRead = 1024, nReadBytes = 0;
unsigned char *pBuffer = new unsigned char[nBytesToRead];
HANDLE hFile;
int Result = ReadFile(hFile, pBuffer, nBytesToRead, &nReadBytes, NULL);
if (!Result)
   cout << "Data not read, error " << GetLastError() << endl;</pre>
else if (nReadBytes != nBytesToRead)
  cout << "Only " << nReadBytes << " bytes instead of " << nBytesToRead << " was read"
       << endl;
else
  cout << nReadBytes << " bytes was read" << endl;</pre>
The complete specification of function ReadFile is on website
```

https://docs.microsoft.com/en-us/windows/desktop/api/fileapi/nf-fileapi-readfile

Asynchronous I/O in Windows (4)

```
To write into file synchronously:
int result name = WriteFile(handle name, pointer to buffer containing data,
number of bytes to write, pointer to variable for number of bytes actually written,
NULL);
If the writing failed, the return value is 0 and GetLastError() returns the error code. The
actual number of written bytes may be less than the number of bytes we wanted to write.
Example:
unsigned long nBytesToWrite = 1024, nWrittenBytes = 0;
unsigned char *pBuffer = new unsigned char[nBytesToWrite];
HANDLE hFile;
int Result = WriteFile(hFile, pBuffer, nBytesToWrite, &nWrittenBytes, NULL);
if (!Result)
   cout << "Data not written, error " << GetLastError() << endl;
else if (nBytesToWrite != nWrittenBytes)
  cout << "Only " << nWrittenBytes << " bytes instead of " << nBytesToWrite
       << " was written" << endl;
else
  cout << nWrittenBytes << " bytes was written" << endl;</pre>
The complete specification of function WriteFile is on website
https://docs.microsoft.com/en-us/windows/desktop/api/fileapi/nf-fileapi-writefile
```

Asynchronous I/O in Windows (5)

The synchronous I/O operations with disk files in most cases do not lead to problems. For Windows, however, any data that is not in the memory region set for the current application is considered to be in a file. For example, if our application controls an external device connected to computer through serial port COM1, then we need to read data from and write data to file \\\.\\COM1\.\\COM1\.\\COM1\.\\COM1\.\\\COM1\) This external device may be temporarily or perpetually (is switched off?) not able to send data to us or retrieve data that we want to send to it. In that case the synchronous I/O operations block their thread and it is impossible to unblock it – in other words the application hangs.

For asynchronous I/O we need Windows events:

```
HANDLE handle_name = CreateEventA (
    attributes, // outside the course scope, set to NULL
    reset_mode, // manual reset TRUE, auto reset FALSE
    initial_state, // not signaled FALSE, signaled TRUE
    name); // outside the course scope, set to NULL
```

Other standard functions for events:

```
SetEvent(handle_name); // state to signaled
ResetEvent(handle_name); // state to non-signaled
CloseHandle(handle_name);
```

See more on website:

https://docs.microsoft.com/en-us/windows/desktop/Sync/synchronization-functions#event-functions

Asynchronous I/O in Windows (6)

The events are to block and unblock a thread:

```
int result name = WaitForSingleObject(event handle name, timeout in ms);
```

If the event is non-signaled, *WaitForSingleObject* function does not return and thus blocks its thread. To unblock, another thread or an asynchronous I/O operation must set the event to signaled. The function returns also when the timeout interval elapses. If the timeout is set to *INFINITE*, time is not measured.

The return value may be:

- WAIT_OBJECT_0 the event was set to signaled
- WAIT_TIMEOUT the event is still non-signaled, but time has elapsed
- *WAIT_FAILED* system problems

If the event is defined as auto-reset, *WaitForSingleObject* function also turns it back to not signaled. If the event is defined as manual reset, the user has to apply function *ResetEvent*.

The other blocking / unblocking function is

It checks the state of several events stored into array. If the third parameter all_or_one is TRUE, function WaitForMultipleObjects returns only when all the specified events are signaled. If all_or_one is FALSE, it returns when at least one of them is signaled. If just the array's i-th event has become signaled, the return value is $WAIT_OBJECT_0 + i$. If several or all the events are signaled, i is the smallest index from the signaled events subset.

Asynchronous I/O in Windows (7)

For asynchronous I/O we need a struct of type *OVERLAPPED*: OVERLAPPED Overlapped;

memset(&Overlapped, 0, sizeof Overlapped); // fill with zeroes

Next we have to create an initially non-signaled event and store its handle in *Overlapped*: Overlapped.hEvent = CreateEventA(NULL, FALSE, FALSE, NULL);

We need at least one more non-signaled event that is triggered from another thread. For example, let us suppose that this event with handle *hExitEvent* will be set to signaled when the user has decided to exit the application.

To finish the preparations create an array of events:

HANDLE hEvents[] = { Overlapped.hEvent, hExitEvent };

In calls to *ReadFile* and *WriteFile* functions the last parameter must be the pointer to *Overlapped*, for example:

unsigned long nBytesToRead = 1024, nReadBytes = 0; unsigned char *pBuffer = new unsigned char[nBytesToRead]; HANDLE hFile;

int Result = ReadFile(hFile, pBuffer, nBytesToRead, &nReadBytes, &Overlapped);

If the output value is not *FALSE*, the I/O operation succeeded without any delays and we may check does the number of read or written bytes match the number of bytes we wanted to read or write (see slides *Asynchronous I/O in Windows (3)* and *(4)*).

Asynchronous I/O in Windows (8)

If the output value is *FALSE*, we must call *GetLastError()*: int error = GetLastError();

If the error code is *ERROR_IO_PENDING*, the I/O operation failed now but may succeed later and we have to wait:

int WaitResult = WaitForMultipleObjects(2, hEvents, FALSE, timeout);

The waiting stops when:

- At last the I/O operation was completed, WaitResult is WAIT_OBJECT_0.
- The user wants to exit the application and triggers hExitEvent, WaitResult is $WAIT_OBJECT_0 + 1$. Thus the application is always under the user's control.
- Timeout occurred (only if *timeout != INFINITE*), *WaitResult* is *WAIT_TIMEOUT*.
- Some system errors, WaitResult has other values

If the I/O operation was completed, we may get the number of bytes that was actually read or written:

GetOverlappedResult(hFile, &Overlapped, &nReadBytes, FALSE); or

GetOverlappedResult(hFile, &Overlapped, &nWrittenBytes, FALSE);

If the error code was not IO_ERROR_PENDING, the I/O operation has totally failed.

See also WindowsAsyncIO.cpp from IAX0587 Examples.zip.

Dynamic link libraries (1)

Large application consist of more that one file: *.exe + several *.dll. Why the DLLs are necessary:

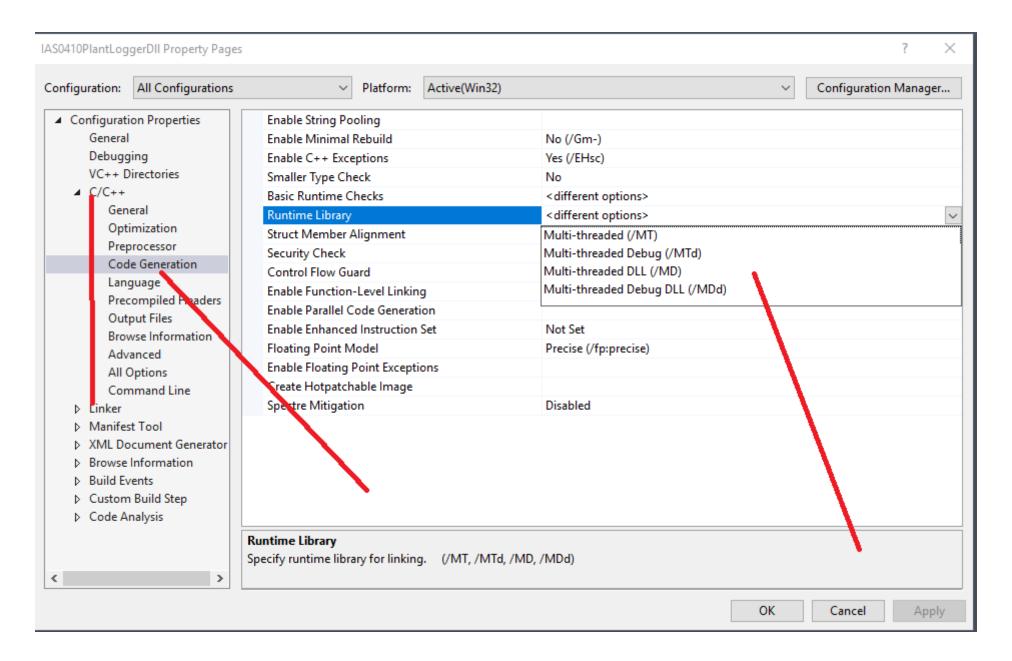
- 1. Very large executive files can be divided into smaller modules.
- 2. In development: each programmer (or group of programmers or company) can link his part of code (a software module) as a DLL and thus work without disturbing the other participants.
- 3. In maintenance: the developer changes one of the DLLs and sends it to the customer.
- 4. Industrial software development: one DLLs may be used in many different applications.

There are two options for connecting *.exe and *.dll:

- 1. Implicit linking: the DLLs are connected to the application when Windows is loading the application into memory.
- 2. Explicit linking: a DLL is connected only when the application needs and calls the *LoadLibrary()* function.

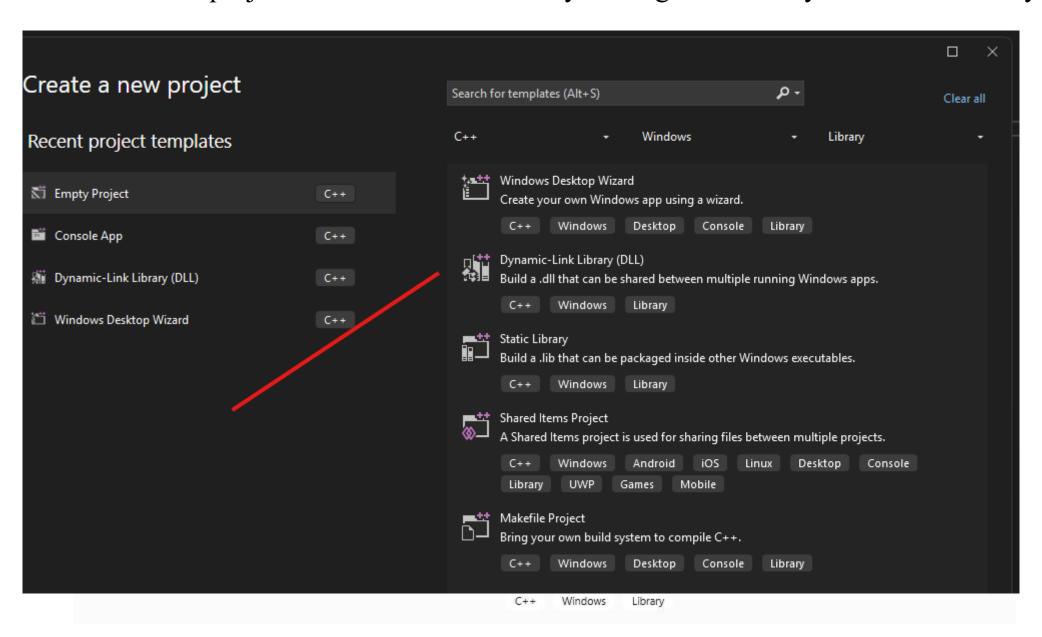
Important for Visual Studio users: Project properties \rightarrow C/C++ \rightarrow Code generation \rightarrow Runtime library has two options: multithreaded /MT (runtime support libraries are linked to the application, the total amount of *.exe is large), or multithreaded DLL/MD (runtime support libraries are applied as DLLs, the amount of *.exe is smaller but when the customer's PC does not have all the necessary libraries, the application crashes).

Dynamic link libraries (2)



Dynamic link libraries (3)

To start with DLL project inform the wizard that your target is C++ Dynamic-link Library.



Dynamic link libraries (4)

Suppose that our project name is *Example*. The wizard creates several files, among them file *dllmain.cpp* (analyzed on the next slide). First create a *.*cpp* file and a *.*h* file (by tradition their names must match the DLL name, so in our case *Example.cpp* and *Example.h*). In the header file write code:

```
#ifdef EXAMPLE_EXPORTS

#define LIBSPEC _declspec(dllexport)

#else

#define LIBSPEC _declspec(dllimport)

#endif
```

Constant *EXAMPLE_EXPORTS* (generally dllname_EXPORTS) is created by wizard and used when the DLL code is compiled. It is not created for applications that use DLLs. All the prototypes and definitions of functions that the DLL exports (i.e. they will be called by applications using this DLL) must start with LIBSPEC, for example in *Example.h*:

```
LIBSPEC int Sum(int, int);
in Example.cpp:
#include "pch.h" // if the wizard created this file
#include "Example.h"
LIBSPEC int Sum(int x1, int x2)
{
    return x1 + x2;
}
```

Dynamic link libraries (5)

```
The wizard-created function DllMain() from dllmain.cpp is the entry point function.
#include "pch.h"
BOOLAPIENTRY DllMain(HMODULE hModule, DWORD ul reason for call,
                            LPVOID lpReserved)
{// DllMain() is called automatically each time when the process or one of its threads
 // attaches or detaches the DLL:
 // 1. In case of implicit linking when the process starts and terminates
 // 2. In case of explicit linking when the application calls LoadLibrary() and
 // FreeLibrary() functions
  switch (ul reason for call)
  case DLL PROCESS ATTACH:
  case DLL THREAD ATTACH:
    // TODO: add code for initialization operations
    break;
  case DLL THREAD DETACH:
  case DLL PROCESS DETACH:
    // TODO: add code for clean-up operations
    break;
  return TRUE;
```

Dynamic link libraries (6)

Visual Studio builds two files: *Example.dll* and *Example.lib*. If the application using our DLL applies explicit linking, it needs only *Example.dll*. In case of implicit linking it needs 3 files: *Example.dll*, *Example.lib* and *Example.h*.

In application with implicit linking calls to functions exported by DLL do not differ from calls to standard functions. The prototypes during compiling are read from *Example.h* that must be copied into the folder where the other source code files are located. However, to link the application we need the object modules (*.obj) of DLL functions. To play a trick on linker we need *Example.lib* that contains stubs – short empty functions replacing the actual DLL functions. When the application is running, instead of stubs the functions from DLL are called. You may copy the *.lib into folder containing source codes.

The easest way is to put DLL into the folder where the *.exe is located.

In case of explicit linking the application must load the DLL:

HMODULE handle_name LoadLibraryA(dll_filename_as_C_string);

Zero handle value means that the DLL was not loaded. To know the reason call *GetLastError()*. To detach the DLL:

FreeLibrary(handle_name);

To call a function we need pointer to it:

FARPROC pointer_name = GetProcAddress(handle_name, function_name_as_C_string);

Zero result value means that the function was not found. To know the reason call *GetLastError()*.

Dynamic link libraries (7)

Explicit linking example: #include "Windows.h" HMODULE hDLL = LoadLibraryA("DLLExample.dll"); if (!hDLL) cout << "DLLExample.dll not found, error " << GetLastError() << endl; return; FARPROC pSum = GetProcAddress(hDLL, "Sum"); if (pSum == NULL) FreeLibrary(hDLL); cout << "Function Sum() not found, error " << GetLastError() << endl;</pre> return; int result = ((int(*)(int, int))pSum)(5, 6); // cast pointer to actual type FreeLibrary(hDLL);

return;

Dynamic link libraries (8)

Turn attention that C++ compiler performs so called name mangling. For example, a function with prototype

```
unsigned char *Run(unsigned char *);
```

in DLL is named as ?Run@@YAPEAEPEAE@Z. To see the new names you can analyse the DLL with Dependency Walker standard utility.

The C compiler keeps the original names. To force the C++ compilers to follow C naming conventions use extern "C" linking:

```
extern "C"
{
    LIBSPEC int Sum(int, int);
}
extern "C"
{
    LIBSPEC int Sum(int x1, int x2)
    {
      return x1 + x2;
    }
}
```