

## Setting up for the lab

1. Restart the SQL Server service (clears stats)
2. Restore your StackOverflow database
3. Copy & run the setup script, will take 1-2 minutes:  
[BrentOzar.com/go/serverlab3](http://BrentOzar.com/go/serverlab3)
4. Start SQLQueryStress:
  1. File Explorer, D:\Labs, run SQLQueryStress.exe
  2. Click File, Open, D:\Labs\ServerLab3.json
  3. Click Go



2.1 p1



**BRENT OZAR**  
UNLIMITED®

**Normal Parallelism**

**CXCONSUMER, CXPACKET, and LATCH\_EX**

2.1 p2

## We covered CPU & storage before this.

Queries go parallel because they:

- Need to read a lot of 8KB pages (PAGEIOLATCH)
- Need to do a lot of CPU work (SOS\_SCHEDULER\_YIELD)

So generally, when I see CXPACKET/CXCONSUMER as the top wait:

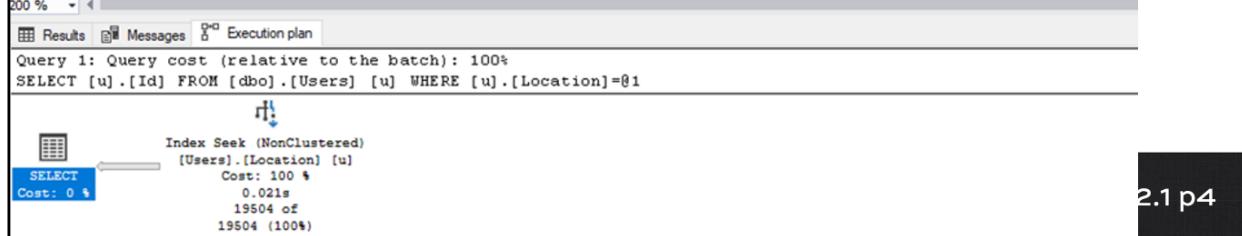
1. Make sure CTFP and MAXDOP have sane defaults  
(and we'll talk about those)
2. Look past CX\*, and look at the next top wait type  
(typically SOS or PAGEIOLATCH)
3. Tune that, and the parallelism usually falls away
4. But sometimes it doesn't (and we cover that in MQT)



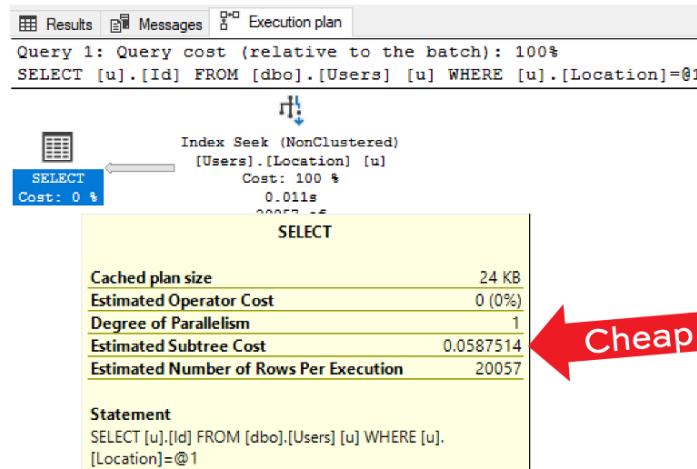
2.1 p3

## Easy query plans only use 1 CPU core.

```
28  CREATE INDEX Location ON dbo.Users(Location);
29  GO
30
31  /* If I need to find all the users in one location: */
32  SELECT u.Id
33  FROM dbo.Users u
34  WHERE u.Location = N'London, United Kingdom';
35  GO
```



Because their cost is low.

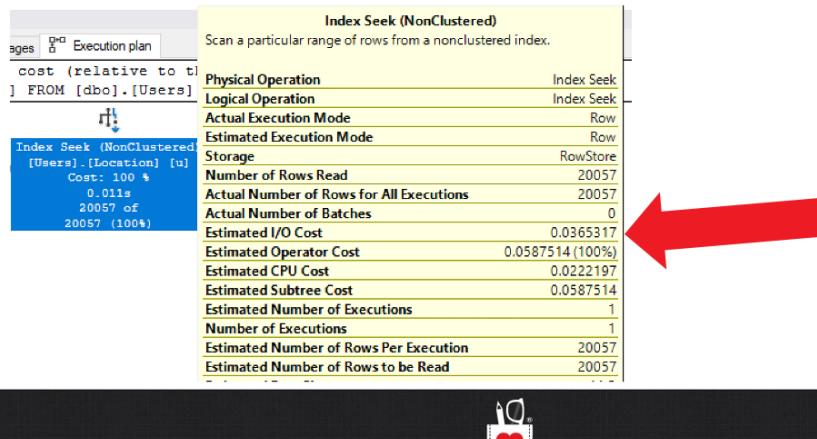


Cheap

2.1 p5

## Each operator has a CPU and IO cost

And every operator's cost is added together to get the Estimated Subtree Cost of the entire query.



Index Seek (NonClustered)	
Scan a particular range of rows from a nonclustered index.	
Physical Operation	Index Seek
Logical Operation	Index Seek
Actual Execution Mode	Row
Estimated Execution Mode	Row
Storage	RowStore
Number of Rows Read	20057
Actual Number of Rows for All Executions	20057
Actual Number of Batches	0
Estimated I/O Cost	0.0365317
Estimated Operator Cost	0.0587514 (100%)
Estimated CPU Cost	0.0222197
Estimated Subtree Cost	0.0587514
Estimated Number of Executions	1
Number of Executions	1
Estimated Number of Rows Per Execution	20057
Estimated Number of Rows to be Read	20057

2.1 p6

## How query costs were calculated: BrentOzar.com/go/nick

I want to share with you a story I heard from Lubor Kollar who was the Program Manager for the Query Optimizer Team of SQL Server ,and now he is leading the SQL Server Customer Advisory Team .(<http://blogs.msdn.com/sqlcat/default.aspx>) which is part of the product group

the story goes that when the new query optimizer was developed for SQL server 7.0 in the Query Optimizer team there was a programmer called nick (I am sorry but I do not know his last name) ,he was responsible for calculating query costs (among other things...) ,and he had to decide how to generate the cost number for a query ,so he decided that if query runs for 1 second on his own pc the cost will be .... 1 ,so we can finally answer the question what is "estimated subtree cost = 1" means,it means ladies and gentleman that it runs for 1 second on nick's machine

: sometime ago Lubor Kollar sent me a picture of nick's machine ,and here it is



.....rumor's are saying that this pc is now stored in Lubor Kollar's garage

2.1 p7

**Kendra Little coined the term “Query Bucks”**

[BrentOzar.com/go/querybucks](https://BrentOzar.com/go/querybucks)



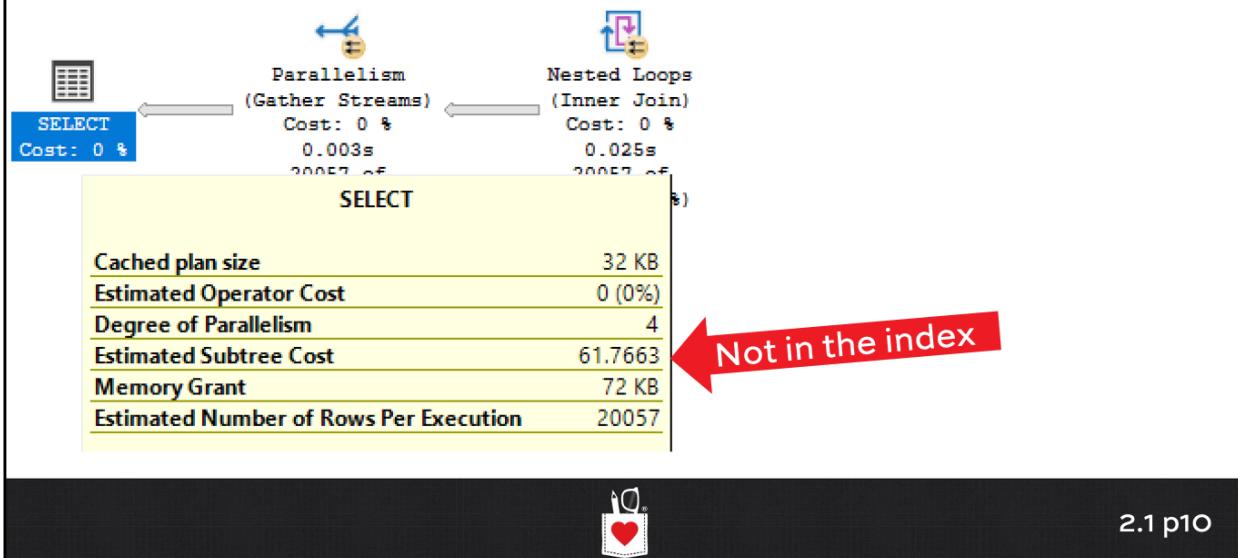
2.1 p8



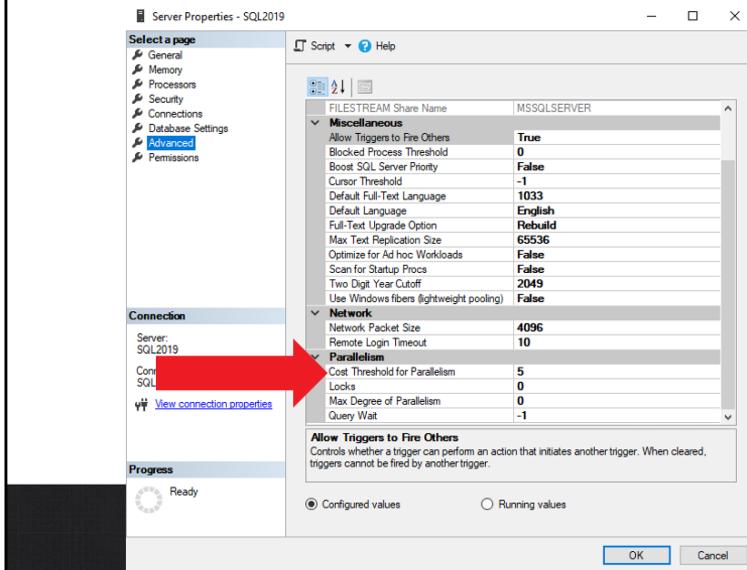
## Add a key lookup, it goes parallel...



Because the query cost went up.



## SQL Server compares the query's total cost to



Server-level setting:  
Cost Threshold for Parallelism

Means queries must be this high  
to ride the parallelism train

2.1 p11

**And as long as our query doesn't have these:**

- Scalar/MSTVFs/Table Variables
  - Computed columns, check constraints
- TOP
- Global aggregates
- Sequence operators
- Backwards scans
- Recursive CTEs
- Some system tables and functions
- CLR that performs data access
- Dynamic cursors
- OUTPUT clause

**Far and away the worst,  
additional overhead**

**Can cause serial zones in plans.**

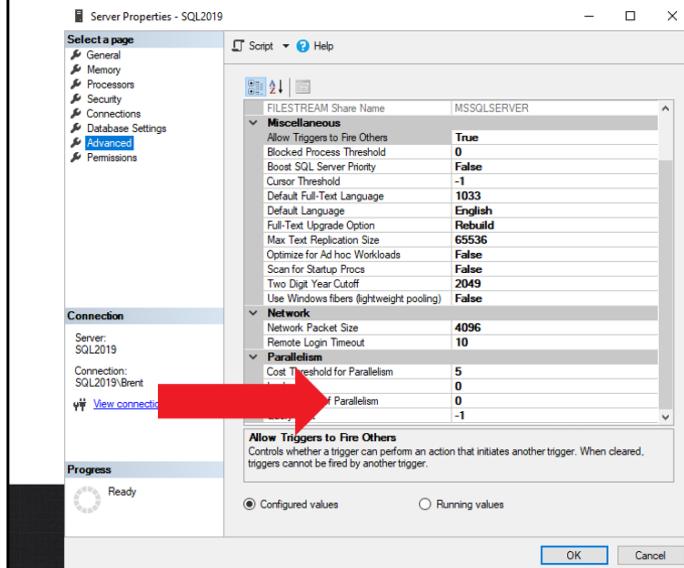
**Far less common. System  
functions and views are less  
likely to cause overall query  
performance issues.**

**Also setting MAXDOP to 1**



2.1 p12

Then consider parallelizing it at MAXDOP cores.



Up to the number of cores  
set by Max Degree of  
Parallelism.

O = unlimited, as in all of the  
server's CPU cores.

(There's some other math  
involved here, but I'm  
keeping it simple.)

2.1 p13

## What MAXDOP controls

MAXDOP controls the maximum number of

- Cores a query can use
- Threads that a parallel branch can use
- Operators in a branch may share threads

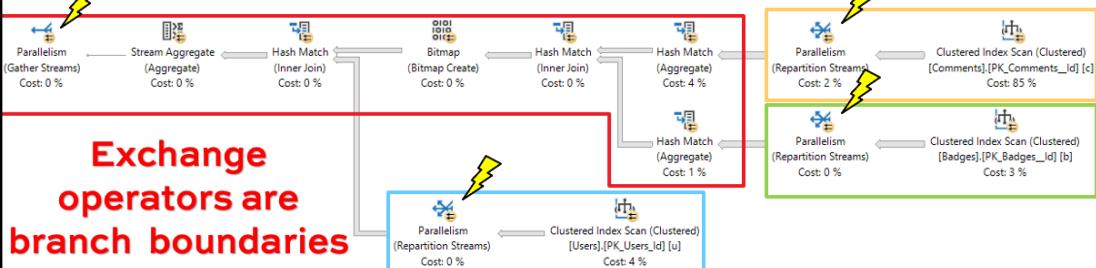
This setting does not control

- The total number of threads a parallel query can use
- The number of threads per parallel operator



2.1 p14

## What counts as a branch?



Exchange operators are branch boundaries

This is the number of branches that can be active concurrently, not the total branches.

ThreadStat	
Branches	2
UsedThreads	8

Not all branches can execute concurrently because of blocking operators (hashes, sorts)

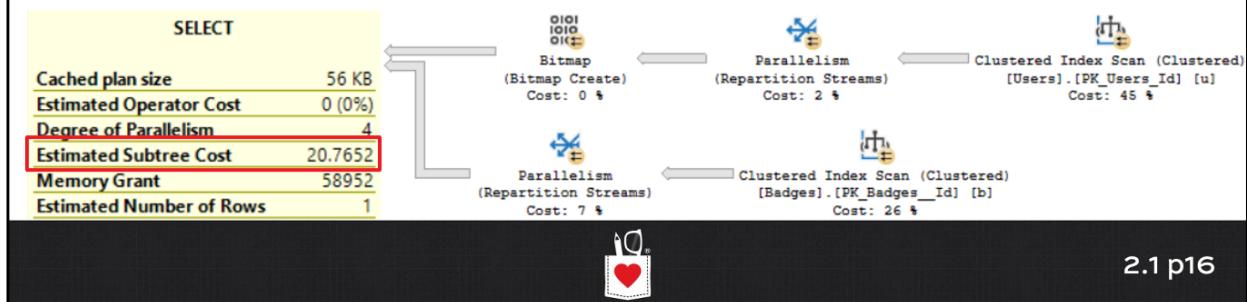


2.1 p15

The parallel plan's cost might be < CTFP

Serial plan cost is what needs to exceed Cost Threshold

To see what the serial plan's cost was,  
run it with an OPTION (MAXDOP 1) hint



## Back to our query. It went parallel, right?

```
38 |  
39 |     SELECT u.Id, u.DisplayName  
40 |     FROM dbo.Users u  
41 |     WHERE u.Location = N'London, United Kingdom';
```

200 % Execution plan

Query 1: Query cost (relative to the batch): 100%  
SELECT [u].[Id],[u].[DisplayName] FROM [dbo].[Users] AS [u] WHERE [u].[Location]=@1  
Missing Index (Impact 99.9%) [Missing Index, sysname,>] ON [d]

**Looks like it did**

2.1 p17

Well, examine the index seek's properties...

The screenshot shows the SQL Server Management Studio interface with an execution plan tree. A context menu is open over a node labeled 'Index Seek (NonClustered)'. The menu has the following structure:

- Logical Operation
- Actual Execution Mode
- Estimated Execution Mode
- ... (ellipsis)
- Save Execution Plan As...
- Show Execution Plan XML...
- Compare Showplan
- Analyze Actual Execution Plan
- Find Node
- Missing Index Details...
- Zoom In
- Zoom Out
- Custom Zoom...
- Zoom to Fit
- Properties

The 'Properties' option at the bottom of the menu is highlighted with a yellow background.

2.1 p18

Properties	
Index Seek (NonClustered)	
Misc	
Actual Execution Mode	Row
Actual I/O Statistics	
Actual Lob Logical Reads	0
Actual Lob Physical Reads	0
Actual Lob Read Aheads	0
Actual Logical Reads	161
Thread 0	8
Thread 1	20
Thread 10	0
Thread 11	0
Thread 12	0
Thread 13	0
Thread 14	0
Thread 15	0
Thread 16	0
Thread 17	0
Thread 18	0
Thread 19	0
Thread 2	0
Thread 20	0
Thread 21	0
Thread 22	0
Thread 23	0
Thread 24	0
Thread 25	0
Thread 26	0
Thread 27	0
Thread 28	0
Thread 29	0
Thread 3	0
Thread 30	0
Thread 31	0
Thread 32	0
Thread 33	0
Thread 34	0
Thread 35	0
Thread 36	0
Thread 37	0
Thread 38	0

And it's not really parallel here.

**Thread 0 = coordinating thread**

- Reads a little to figure out how to pass out work to the other threads
- He's not actually handling rows though

**Thread 1 = doing some work**

**The rest = bored out of their gourds**



2.1 p19

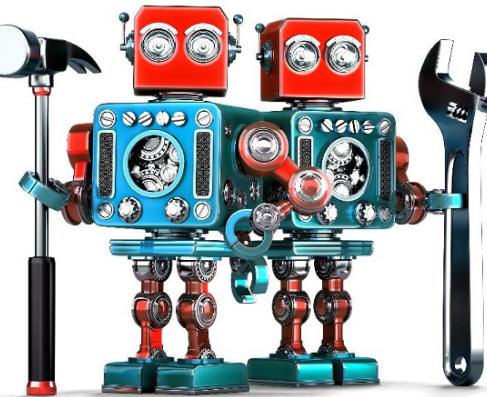
## **When the robots are working right**

**Work is divided across threads**

**Each thread performs equal work**

**Seamless results returned to user**

**Each thread needs – and gets –  
exactly the same amount of  
memory and CPU time**



2.1 p2O

**There's a Robot Overlord who tells the other robots to go get rows**

**There are Robot Workers who obey the Robot Overlord**



2.1 p21

Properties	
Clustered Index Scan (Clustered)	
<input type="checkbox"/>	Misc
Actual Execution Mode	Row
<input type="checkbox"/>	Actual I/O Statistics
<input type="checkbox"/>	Actual Lob Logical Reads
<input type="checkbox"/>	Actual Lob Physical Reads
<input type="checkbox"/>	Actual Lob Read Aheads
<input type="checkbox"/>	Actual Logical Reads
Thread 0	17
Thread 1	48874
Thread 10	43973
Thread 11	48088
Thread 12	47664
Thread 13	45136
Thread 14	45296
Thread 15	46527
Thread 16	45813
Thread 2	42557
Thread 3	43301
Thread 4	45413
Thread 5	45189
Thread 6	43593
Thread 7	49246
Thread 8	49214
Thread 9	48008
<input type="checkbox"/>	Actual Physical Reads
	0

## Problem 1: dividing the input

Sometimes pages aren't evenly divided.

The threads aren't doing the same amount of work.

They're not sharing CPU time or memory either:  
each one is pinned to a specific core, and given a  
specific amount of memory.

One thread can spill to disk while others don't.

One thread can compete with other tasks running  
on the same physical CPU core.



2.1 p22

## Problem 2: unpredictable output

We simply don't know which threads will find results.

For example, if we don't have an index  
and we have to do a table scan to find matches,  
we just don't know which threads will find matches.

The matches could be anywhere in the table.

One thread might just find all of them.

<input checked="" type="checkbox"/> Actual Number of Batches	0
<input type="checkbox"/> Actual Number of Rows	2
Thread 0	0
Thread 1	0
Thread 10	0
Thread 11	0
Thread 12	0
Thread 13	0
Thread 14	0
Thread 15	0
Thread 16	0
Thread 2	0
Thread 3	1
Thread 4	1
Thread 5	0
Thread 6	0
Thread 7	0
Thread 8	0
Thread 9	0
<input type="checkbox"/> Actual Rebinds	0



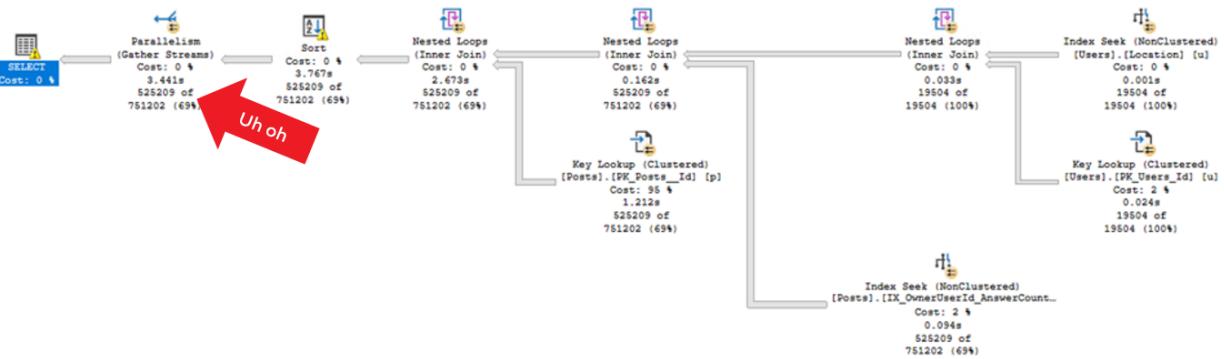
2.1 p23

## When unpredictable output is bad

By itself – just reading out rows – it's not a big deal.

But what if *several* operators are daisy-chained together in the same parallel operation, like we cover in Mastering Query Tuning?

One thread may do all of the work, period.



## Problem 3: odd parallelism bugs

I call these bugs, but Microsoft disagrees:

- Index spools say they're parallel, but they're not
- Order preserving exchanges like parallel merge joins perform terribly

We'll cover these in a separate module on Abnormal Parallelism.

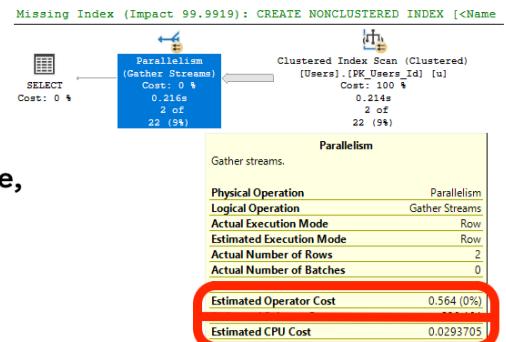


2.1 p25

## Problem 4: parallelism itself has a cost, too.

This one's cheap, but the more data you have, the higher the cost becomes.

Plus, in bigger plans, we have to rebalance the work as we move through the operators, distributing streams of data differently to different threads to compensate for the fact that we don't know which threads will find work.



2.1 p26

## How we measure parallelism

Normally, we use wait stats to measure stuff.

Parallelism makes things tricky though:

- Every robot tracks waits
- The overlord tracks waits too

While a parallel query is running,  
the robot overlord (thread 0) piles up waits  
even though he's not really waiting on things,  
and we can't possibly get rid of his "wait" time.

Robot workers also track "wait" time when they're done  
with their task, and "waiting" on other robots to finish.



2.1 p27

## This is where it gets really terrible.

The robot overlord will always be waiting, and that's totally okay.

The robot workers should NOT be waiting, because that indicates:

- Problem 1: we didn't balance work evenly
- Problem 2: the workers found rows that caused them to do unbalanced work
- Problem 3: there was too much overhead in passing out the data

But for decades, both the robot overlord's wait time and the robot workers' wait time had the same wait type.



2.1 p28

## The CXPACKET debacle

Microsoft finally fixed this in 2016 SP2 & 2017 CU3:

	Robot Overlord	Robot Worker
SQL Server 2017 CU3	CXCONSUMER ("OK")	CXPACKET
SQL Server 2017 pre-CU3	CXPACKET	CXPACKET
SQL Server 2016 SP2	CXCONSUMER ("OK")	CXPACKET
SQL Server 2016 pre-SP2	CXPACKET	CXPACKET
SQL Server 2014 & prior	CXPACKET	CXPACKET

CXCONSUMER was added to separate the coordinator thread waiting on child threads to finish, vs. child threads waiting on each other to finish



2.1 p29

## But it's still really confusing.

CXCONSUMER *should* be totally ignorable – but it's not.

We've found queries in the wild that *only* produce CXCONSUMER waits during unbalanced parallelism.

CXPACKET *should* mean we have a parallelism problem, but which one?

- Problem 1: we didn't balance work evenly
- Problem 2: the workers found rows that caused them to do unbalanced work
- Problem 3: odd parallelism bugs
- Problem 4: too much overhead passing out the data

You just can't tell from wait stats.



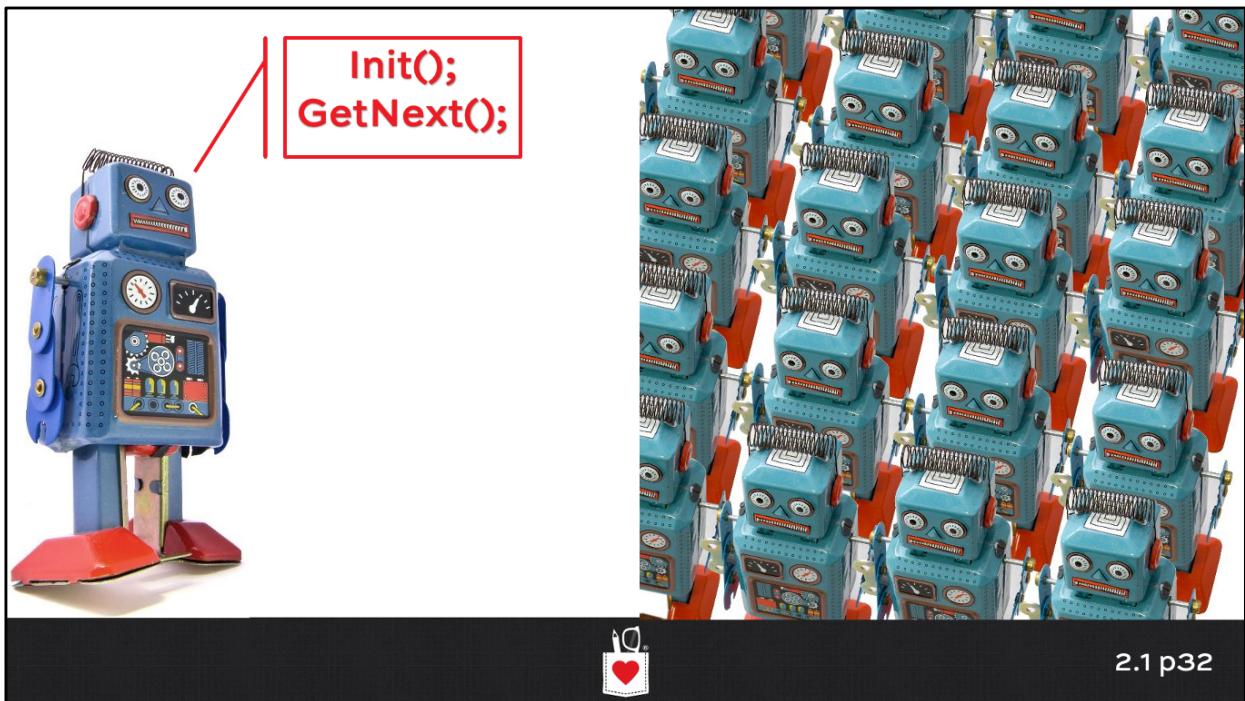
2.1 p30

## You run a “big” query

```
SELECT COUNT(*) AS records
FROM dbo.Posts AS p
JOIN dbo.Votes AS v
    ON v.PostId = p.Id
JOIN dbo.Comments AS c
    ON c.PostId = v.PostId
JOIN dbo.Users AS u
    ON u.Id = p.OwnerUserId
WHERE v.PostId < 10000000
```



2.1 p31

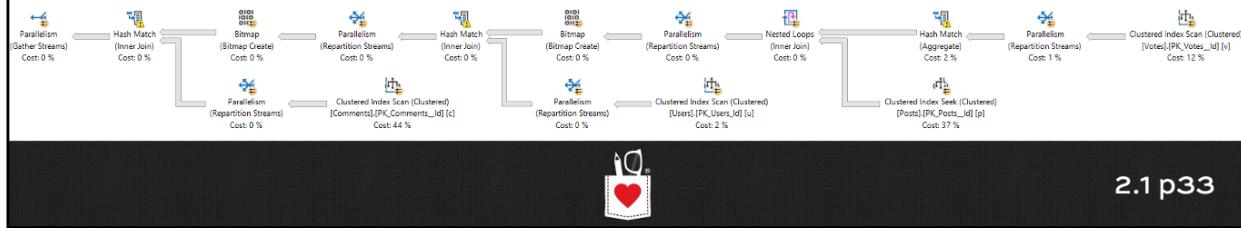


2.1 p32

## This does a lot of CPU work for ~25 seconds

```
EXEC dbo.sp_BlitzFirst @Seconds = 30, @ExpertMode = 1;
```

Seconds Sample	wait_type	wait_category	Wait Time (Seconds)	Per Core Per Second	Signal Wait Time (Seconds)	Percent Signal Waits	Number of Waits	Avg ms Per Wait
30	CXPACKET	Parallelism	63.2	0.2	29.9	47.3	17861	3.5
30	CXCONSUMER	Parallelism	45.3	0.1	26.3	58.1	24440	1.9
30	SOS_SCHEDULER_YIELD	CPU	16.6	0.0	16.6	100.0	16391	1.0
30	PREEMPTIVE_OS_WRITETFILE	Preemptive	0.1	0.0	0.0	0.0	4	19.0
30	RESERVED_MEMORY_ALLOCATION_EXT	Memory	0.1	0.0	0.0	0.0	208370	0.0
30	SOS_PHYS_PAGE_CACHE	Other	0.0	0.0	0.0	0.0	1	1.0
30	LATCH_EX	Latch	0.0	0.0	0.0	0.0	12	0.1
30	MEMORY_ALLOCATION_EXT	Memory	0.0	0.0	0.0	0.0	4539	0.0



2.1 p33

## **But is parallelism really the problem?**

If you looked at a server with queries like this running all the time, you'd see a whole lot of parallelism waits

We see a lot of servers like this

- They've been up for 10 days
- With 20-30 days of parallelism waits

If we change MAXDOP to 1,  
will performance get better?



2.1 p34

## All the single threadies

```
SELECT COUNT(*) AS records
FROM dbo.Posts AS p
JOIN dbo.Votes AS v
    ON v.PostId = p.Id
JOIN dbo.Comments AS c
    ON c.PostId = v.PostId
JOIN dbo.Users AS u
    ON u.Id = p.OwnerUserId
WHERE v.PostId < 10000000
OPTION(MAXDOP 1)
```



2.1 p35

**Not exactly**

**EXEC dbo.sp\_BlitzFirst @Seconds = 90, @ExpertMode = 1;**

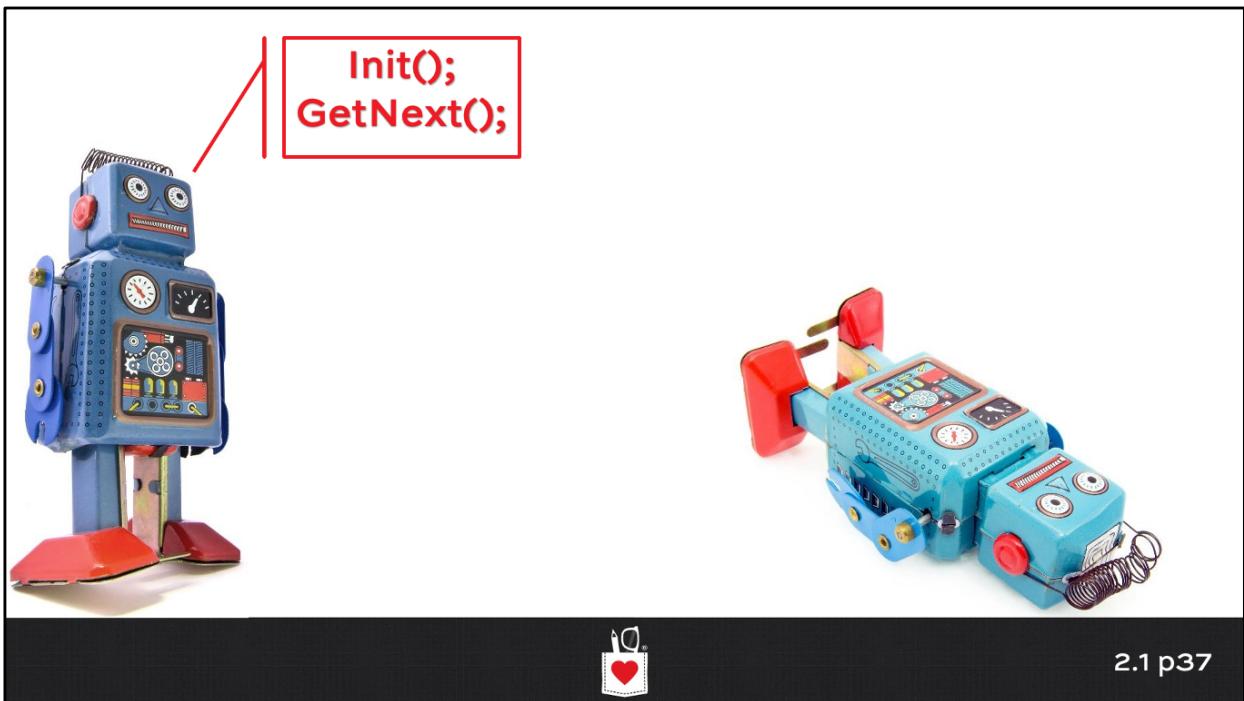
Seconds Sample	wait_type	wait_category	Wait Time (Seconds)	Per Core Per Second	Signal Wait Time (Seconds)	Percent Signal Waits	Number of Waits	Avg ms Per Wait
90	PREEMPTIVE_OS_WRITEFILE	Preemptive	0.1	0.0	0.0	0.0	4	15.5
90	RESERVED_MEMORY_ALLOCATION_EXT	Memory	0.1	0.0	0.0	0.0	201816	0.0
90	SOS_SCHEDULER_YIELD	CPU	0.0	0.0	0.0	0.0	21250	0.0
90	MEMORY_ALLOCATION_EXT	Memory	0.0	0.0	0.0	0.0	7524	0.0

**This runs for  
nearly 90 seconds    It got 3x slower,  
but we got rid of  
parallelism**

**Yay?**



2.1 p36



## So when you see CXPACKET/CXCONSUMER

1. Set CTFP & MAXDOP per industry best practices.
2. Let those settings bake in for a few days.
3. Make index improvements that can reduce times when queries go parallel to scan lots of data & sort it.
4. Let those settings bake in for a few days.
5. Look for queries still doing a lot of work, then tune them.
6. Let those bake in for a few days.
7. Finally, if parallelism is still a serious problem (like 10X wait time ratio), start looking for specific queries with parallelism issues.



2.1 p38

## **1. Setting Cost Threshold for Parallelism**

**Default: 5 query bucks**

**I default to: 50 query bucks**

**I'm fine with a range of: 30-100**

**But be wary: the higher you set it,  
remember that multi-threaded queries  
may suddenly start going single-threaded.**

**Changing this setting clears the plan cache.**



2.1 p39

## 1. Setting Max Degree of Parallelism

Microsoft says: [https://support.Microsoft.com/kb/2806535](https://support.microsoft.com/kb/2806535)

### SQL Server 2016 (13.x) and higher

Starting with SQL Server 2016 (13.x), during service startup if the Database Engine detects more than eight physical cores per NUMA node or socket at startup, soft-NUMA nodes are created automatically by default. The Database Engine takes care of placing logical processors from the same physical core into different soft-NUMA nodes. The recommendations in the table below are aimed at keeping all the worker threads of a parallel query within the same soft-NUMA node. This will improve the performance of the queries and distribution of worker threads across the NUMA nodes for the workload.

Starting with SQL Server 2016 (13.x), use the following guidelines when you configure the **max degree of parallelism** server configuration value:

Server with single NUMA node	Less than or equal to 8 logical processors	Keep MAXDOP at or below # of logical processors
Server with single NUMA node	Greater than 8 logical processors	Keep MAXDOP at 8
Server with multiple NUMA nodes	Less than or equal to 16 logical processors per NUMA node	Keep MAXDOP at or below # of logical processors per NUMA node
Server with multiple NUMA nodes	Greater than 16 logical processors per NUMA node	Keep MAXDOP at half the number of logical processors per NUMA node with a MAX value of 16



2.1 p40

### 3. Making index improvements to reduce CX%

After letting the parallelism changes bake in for a few days:

`sp_BlitzIndex @GetAllDatabases = 1`

- Look for high-value missing indexes
- Remember that they're straight from Clippy, though:  
use the Mastering Index Tuning techniques

`sp_BlitzCache @SortOrder = 'reads' and = 'cpu'`

- Do manual index tuning to reduce scans, sorts
- Optionally, while you're there: tune the queries

Then let these bake in for a few days too.



2.1 p41

## 5. Look for queries still doing a lot of work

After letting the index changes bake in for a few days:

```
sp_BlitzCache @SortOrder = 'reads' and = 'cpu'
```

- Do manual index tuning to reduce scans, sorts
- Optionally, while you're there: tune the queries

And most of the time, that process just obliterates CX% waits.



2.1 p42

## **But sometimes, 1-6 aren't enough because specific queries have abnormal parallelism.**

1. Set CTFP & MAXDOP per industry best practices.
2. Let those settings bake in for a few days.
3. Make index improvements that can reduce times when queries go parallel to scan lots of data & sort it.
4. Let those settings bake in for a few days.
5. Look for queries still doing a lot of work, then tune them.
6. Let those bake in for a few days.
7. Finally, if parallelism is still a serious problem (like 10X wait time ratio), start looking for specific queries with parallelism issues.



2.1 p43

# Recap



2.1 p44

## Start with the basics.

Queries go parallel because they:

- Need to read a lot of 8KB pages (PAGEIOLATCH, LATCH\_EX)
- Need to do a lot of CPU work (SOS\_SCHEDULER\_YIELD)

So when I see CXPACKET/CXCONSUMER as the top waits:

1. Make sure CTFP and MAXDOP have sane defaults (CTFP 30-100, MAXDOP per KB #2806535)
2. Do normal index tuning and query tuning
3. The parallelism usually falls away – but if it doesn't, review the abnormal parallelism causes & fixes in Mastering Query Tuning



2.1 p45