



**IBM zEnterprise EC12**  
**CPU Facilities**

Dan Greiner (dgreiner@us.ibm.com)  
IBM z/Server Architecture

SHARE 120, Session 12670  
Tuesday, 5 February 2013, 3:00 p.m.

IBM Systems and Technology Group (STG)

© Copyright International Business Machines Corporation 2012, 2013

This presentation will discuss the major new CPU facilities added to the IBM zEnterprise EC12 system.

Earlier versions of this presentation were provided to members of the IBM Early-Support Program for the zEC12 in the summer of 2012 and at the IBM Technical Disclosure Meeting in the autumn of 2012. If you attended such a presentation, please be advised that the SHARE-120 version contains numerous corrections and clarifications.



## The Legal Stuff

- The following terms are registered trademarks of the International Business Machines Corporation in the United States, other countries, or both:
  - ▶ IBM
  - ▶ IBM logo
- The following terms are trademarks of the International Business Machines Corporation in the United States, other countries, or both:
  - ▶ ESA/390
  - ▶ z/Architecture
  - ▶ z/OS
  - ▶ z/VM
- The following are trademarks or registered trademarks of other companies:
  - ▶ IEEE is a trademark of the Institute of Electrical and Electronics Engineers, Inc. in the United States, other countries, or both.
  - ▶ Java is a trademark of Oracle America, Inc. in the United States, other countries, or both.
  - ▶ Linux is a registered trademark of Linus Torvalds in the United States, other countries or both.
  - ▶ Unicode is a registered trademark of Unicode, Incorporated in the United States, other countries, or both.
  - ▶ Other trademarks and registered trademarks are the properties of their respective companies.
- All information contained in this document is subject to change without notice. The products described in this document are not intended for use in applications such as implantation, life support, or other hazardous uses where malfunction could result in death, bodily injury or catastrophic property damage. The information contained in this document does not affect or change IBM product specifications or warranties. Nothing in this document shall operate as an express or implied license or indemnity under the intellectual property rights of IBM or third parties. All information contained in this document was obtained in specific environments, and is presented as an illustration. The results obtained in other operating environments may vary.
- While the information contained herein is believed to be accurate, such information is preliminary, and should not be relied upon for accuracy or completeness, and no representations or warranties of accuracy or completeness are made.
- The information contained in this document is provided on an "AS IS" basis. In no event will IBM be liable for damages arising directly or indirectly from any use of the information contained in this document.
- This publication was produced in the United States. IBM may not offer the products, services or features discussed in this document in other countries, and the information may be subject to change without notice. Consult your local IBM business contact for information on the product or services available in your area.
- All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only

This slide reviews the trademarks that may be shown in the presentation. Also, this slide includes various disclaimers as to the content of the presentation.

## Topics du Jour

- **Interlocked-access facility 2**
- **DFP zoned-conversion facility**
- **Execution-hint facility**
- **Load-and-trap facility**
- **Miscellaneous general instructions**
- **Transactional-execution facility**
- **Processor-assist facility**
- **Enhanced-DAT facility 2**
- **Local-TLB-clearing facility**

This slide enumerates the new CPU facilities introduced in the IBM z/Enterprise EC12. Each of these facilities will be discussed in detail in the subsequent slides.

The session describes the CPU facilities and instructions in detail, and assumes a familiarity with assembler language, machine-instruction formats, and basic z/Architecture.

Much of the material described in today's presentation is related to the characteristics of multiprocessing ... particularly, in improving the performance of MP applications that share common memory locations.

## Interlocked-Access Facility 2

- **System z196 introduced the interlocked-access facility**
  - ▶ **LOAD & ADD, LOAD & ADD LOGICAL, LOAD & AND, LOAD & OR, LOAD & XOR, LOAD PAIR DISJOINT**
  - ▶ **Now called interlocked-access facility 1**
- **Interlocked-access facility 2 provides:**
  - ▶ **Guaranteed block-concurrent interlocked update for:**
    - **ADD IMMEDIATE (ASI, AGSI) with aligned operands**
    - **ADD LOGICAL IMMEDIATE (ALSI, ALGSI) with aligned operands**
    - **AND (NI, NIY)**
    - **OR (OI, OIY)**
    - **EXCLUSIVE OR (XI, XIY)**
  - ▶ **Facility indication bit 52**
  - ▶ **Facility is retrofit to the z196!**

The interlocked-access facility was introduced in the System z196 processor in September of 2010, and included the instructions listed on this slide.

These instructions provide improved performance for certain sequences of operations that may be executed in a multiprocessing environment. The instructions provide a block-concurrent, interlocked update for loading, performing and operation, and storing a result (commonly known as an atomic operation). This facility is now called the interlocked-access facility 1.

The interlocked-access facility 2 provides an assurance that the instructions listed here will also perform in an interlocked, block-concurrent manner. Many of these instructions such as AND (NI), EXCLUSIVE OR (XI), and OR (OI) have existed since the original S/360, with programming notes advising that they cannot safely be used in an MP environment. The interlocked-access facility 2 changes that, assuring that these instructions will perform in an interlocked manner.

The interlocked-access facility 2 was actually present in the System z196, but no facility indication was originally provided. Through a firmware upgrade, the interlocked-access facility 2 indication is now provided on all z196 processors.

## DFP Zoned-Conversion Facility

- **Adds instructions for converting between DFP and zoned format**
  - ▶ **May provide substantial performance improvement for applications that use packed-decimal data**
  - ▶ **By converting to DFP and performing calculations using DFP instructions, numerous storage accesses may be avoided**
  - ▶ **Four new instructions:**
    - Long / extended DFP format
    - To / from zoned format
  - ▶ **Facility indication bit 48**
- **Formally documents zone code 0011 binary (ASCII format)**

The DFP zoned-conversion facility provides four new instructions for converting between the zoned-format in storage and a decimal-floating-point (DFP) format value in a floating-point register.

Applications that use zoned and/or packed data formats may yield increased performance by adapting to perform the arithmetic operations in DFP, while retaining the legacy packed or zoned formats.

Also, as a part of this update, the architecture has also been adapted to formally document the zone code of 0011 binary as representing the ASCII numeric zone.

## DFP Zoned-Conversion Facility (2)

- **CONVERT FROM ZONED**

CDZT  $R_1, D_2(L_2, B_2), M_3$  [RSL-b] (long DFP result)

ED	L <sub>2</sub>	B <sub>2</sub>	D <sub>2</sub>	R <sub>1</sub>	M <sub>3</sub>	AA
----	----------------	----------------	----------------	----------------	----------------	----

CXZT  $R_1, D_2(L_2, B_2), M_3$  [RSL-b] (extended DFP result)

ED	L <sub>2</sub>	B <sub>2</sub>	D <sub>2</sub>	R <sub>1</sub>	M <sub>3</sub>	AB
----	----------------	----------------	----------------	----------------	----------------	----

- **M<sub>3</sub> field:**

0 – Sign control (S: 0=2<sup>nd</sup> operand unsigned; 1=2<sup>nd</sup> operand signed)

1 – reserved

2 – reserved

3 – reserved

- **CXZT 2<sup>nd</sup>-operand can provide a length of 34 digits!**

The CONVERT FROM ZONED instructions are of the RSL instruction format (subformat b). There are instructions for converting a zoned value that result in either a long (64-bit) DFP value or an extended (128-bit) value.

The first operand in the R<sub>1</sub> field specifies a floating-point register into which the DFP-format result is placed.

The second operand is the address of a storage-operand containing the zoned value to be converted; the L<sub>2</sub> field indicates the length of the zoned value in bytes.

The M<sub>3</sub> field contains a one-bit sign control which indicates whether the second-operand is to be treated as a signed or unsigned value.

Note, because of the capacity of DFP number representations, CXZT is capable of accommodating a 34-digit length ... substantially larger than can be accommodated by normal packed-decimal instructions.

## DFP Zoned-Conversion Facility (3)

### ■ CONVERT TO ZONED

CZDT  $R_1, D_2(L_2, B_2), M_3$  [RSL-b] (long DFP source)

ED	L <sub>2</sub>	B <sub>2</sub>	D <sub>2</sub>	R <sub>1</sub>	M <sub>3</sub>	A8
----	----------------	----------------	----------------	----------------	----------------	----

CZXT  $R_1, D_2(L_2, B_2), M_3$  [RSL-b] (extended DFP source)

ED	L <sub>2</sub>	B <sub>2</sub>	D <sub>2</sub>	R <sub>1</sub>	M <sub>3</sub>	A9
----	----------------	----------------	----------------	----------------	----------------	----

### ■ M<sub>3</sub> field:

- 0 – Sign control (S: 0=2<sup>nd</sup> operand unsigned; 1=2<sup>nd</sup> operand signed)
- 1 – Zone control (Z: 0=zone stored as 1111; 1=zone stored as 0011)
- 2 – Plus-sign-code control (P: 0=plus is 1100 binary; 1=plus is 1111)
- 3 – Force-plus-zero control (F: 0=-0 unchanged; 1=-0 made +)

### ■ CZXT 2<sup>nd</sup>-operand can accommodate a 34-digit result!

The CONVERT TO ZONED instructions are of the RSL instruction format (subformat b). There are instructions for converting either a long or extended DFP value into a zoned value.

The first operand is the R<sub>1</sub> field specifies a floating-point register containing the DFP number to be converted.

The second operand is the address of a storage-operand into which the zoned result will be placed; the L<sub>2</sub> field indicates the length of the zoned value in bytes.

The M<sub>3</sub> field contains four separate controls:

- Bit zero controls the sign of the result.
- Bit one controls the resulting zone (0 means EBCDIC, 1 means ASCII)
- Bit 2 specifies the encoding of a positive sign value in the result
- Bit 3 indicates whether a DFP -0 value should be made positive.

Note, because of the capacity of DFP number representations, CZXT is capable of accommodating a 34-digit length ... substantially larger than can be accommodated by normal packed-decimal instructions.

## Execution-Hint Facility

- Provides the following instructions:
  - ▶ BRANCH PREDICTION PRELOAD
  - ▶ BRANCH PREDICTION RELATIVE PRELOAD
  - ▶ NEXT INSTRUCTION ACCESS INTENT
- When the facility is installed, these instructions provide hints to the CPU as to anticipated branches and operand accesses
  - ▶ May provide performance improvement (if used properly)
  - ▶ May degrade performance (if abused)
  - ▶ Otherwise, instructions act as no-ops, and do not affect conceptual sequence of execution.
- Facility indication bit 49

The execution-hint facility provides three instructions which can be used to provide hints to the CPU as to various branching conditions and the storage-access intent of a subsequent instruction. These instructions are used by IBM compilers to optimize instruction flow in the CPU pipeline.

When properly used, these instructions may improve performance. However, when improperly used, these instructions may actually degrade performance by mis-directing CPU branch prediction logic and prefetching controls.

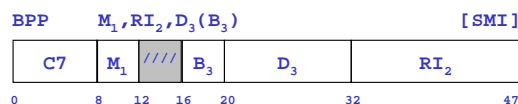
However, regardless of how the instructions are used, they otherwise act as no-operation instructions (no-ops), and do not affect the logic of program execution.

This facility – as well as several others in the zEC12 – are indicated by facility indication 49 (as stored by STORE FACILITY LIST EXTENDED).

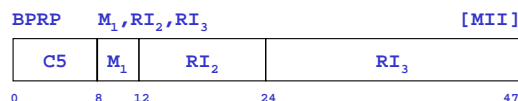


## Branch Prediction (1)

### ■ BRANCH PREDICTION PRELOAD



### ■ BRANCH PREDICTION RELATIVE PRELOAD



### ■ Provides branch hint to CPU

- ▶ M<sub>1</sub> field contains branch-type code (next slide)
- ▶ 2<sup>nd</sup> operand is relative address of a branching instruction
- ▶ 3<sup>rd</sup> operand is address of expected branch location

BRANCH PREDICTION PRELOAD and BRANCH PREDICTION RELATIVE PRELOAD are instructions to provide the CPU's branch-prediction logic with a direct assertion of the programmer's intent for a branch instruction.

In both instructions:

- The M<sub>1</sub> field contains a code designating the type of branch instruction designated by the second operand (see the next slide for details).
- The RI<sub>2</sub> field contains a signed relative-immediate address (relative to the PSW instruction address) that designates the branch instruction.

The third operand designates the anticipated branch location of the instruction designated by the second operand. For BPP, the third operand is a classic base-and-displacement form, and for BPRP, the third operand is a 24-bit (!) signed value that is relative to the PSW instruction address.

## Branch Prediction (2)

### BRANCH PREDICTION PRELOAD $M_1$ Codes

$M_1$ Code	Inst. Leng	Corresponding Branch Instruction (designated by $RI_2$ field)	Usage
0	4	BC	Branch table
1-4	--	Reserved	--
5	2	BALR, BASR, BCR	Static calling linkage
6	2	BCR	Returning linkage
7	2	BALR, BASR, BCR	Dynamic calling linkage
8	4	BC, BCT, BRXH, BRXLE, BXH, BXLE, BRC, BRCT, BRCTG, BCTGR	Cond. or uncond. branches
9	4	BAL, BAS, BRAS	Static calling linkage
10	4	BC	Uncond. return linkage
11	4	BAL, BAS	Dynamic calling linkage
12	6	BRCTH, BRCL, BCTG, BXHG, BXLEG, BRXHG, BRXLG, CGRJ, CLGRJ, CRJ, CLRJ, CGIJ, CLGIJ, CIJ, CLIJ, CGRB, CLGRB, CRB, CLRB, CGIB, CLGIB, CIB, CLIB	Conditional or unconditional branches
13	6	BRASL	Static calling linkage
14	4	EX	
15	6	EXRL	

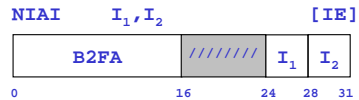
This slide enumerates the code values that may be specified in the  $M_1$  field.

Note that the same instruction appears for different codes. For example, BALR, BASR, and BCR are used in both codes 5 and 7, and BCR also appears in code 6. Code 5 indicates that the designated branch instruction is used for calling a subroutine, and the target location is expected to always be a single location. Code 7 also indicates that the designated branch instruction is used for calling a subroutine, but the target location may be dynamically determined by the program. Code 6 indicates a BCR instruction that is used to return from a called subroutine. Similar abstractions appear for codes 9 and 11.

Performance may be degraded if the second operand does not designate a branch instruction that is used in accordance with the  $M_1$  encoding, or if the branch instruction does not branch to the location specified by the third operand.

## Next-Instruction Access Intent (3)

### ■ NEXT INSTRUCTION ACCESS INTENT



### ■ Provides hint to CPU as to storage use of next-sequential instruction

- ▶  $I_1$  corresponds to the lowest-numbered storage operand (if any)
- ▶  $I_2$  corresponds to the 2<sup>nd</sup>-lowest-numbered storage operand (if any)
- ▶  $I_1$  and  $I_2$  encodings:
  - 0 – Corresponding operand may or may not be accessed
  - 1 – Corresponding operand will be stores, and may be fetched
  - 2 – Corresponding operand will be fetched
  - 3 – Corresponding operand will not be accessed

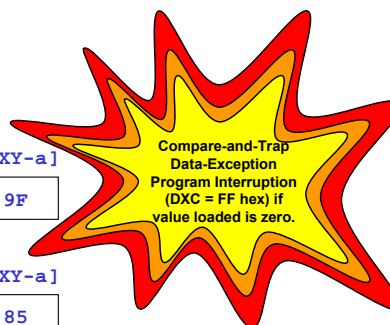
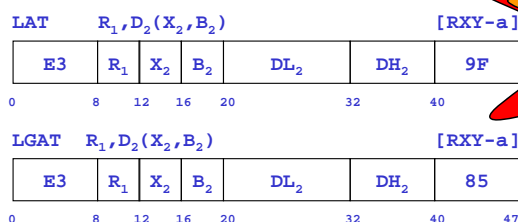
The NEXT INSTRUCTION ACCESS INTENT instruction provides a means for the program to indicate the anticipated use of the storage locations designated by the next instruction.

The two operands are each 4-bit immediate fields that indicate the anticipated usage of the storage operand(s) of the next sequential instruction. The  $I_1$  field represents the lowest-numbered storage operand, and the  $I_2$  field represents the second-lowest-numbered storage operand (if any). The encodings are listed on the slide.

Performance may be degraded if the subsequent instruction uses its storage operands differently than that specified by the respective  $I_1$  and  $I_2$  fields.

## Load-and-Trap Facility (1)

- Provides equivalent function to LOAD, LOAD HIGH, LOAD LOGICAL, and LOAD LOGICAL 31 BITS, but ...
  - ▶ Causes a compare-and-trap data exception if the designated storage operand contains zero
  - ▶ PIC 0007 hex, DXC FF hex
- Facility indication bit 49
- LOAD AND TRAP



The load-and-trap facility provides a means by which a value can be loaded from storage; if the value contains zero, then a compare-and-trap data exception is recognized. The instructions are equivalent to executing a load instruction followed by a compare-and-trap instruction with a comparand of zero.

Each instruction is of the RXY instruction format (subformat –a), meaning that the second operand designates a storage location by means of a base register, an index register, and a 20-bit signed displacement field (that is, a long displacement). The result is loaded into the register designated by the R<sub>1</sub> field.

- LAT loads a 32-bit value into bits 32-63 of the register, and leaves the remaining bits unchanged.
- LGAT loads a 64-bit value into bits 0-63 of the register.

In either case, if the value loaded is zero, then a compare-and-trap data exception is recognized. The program-interruption code is 0007, and the data-exception code is FF hex.

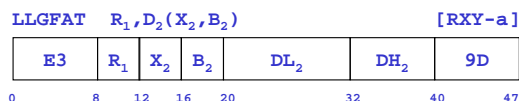
This facility – as well as several others in the zEC12 – are indicated by facility indication 49 (as stored by STORE FACILITY LIST EXTENDED).

## Load-and-Trap Facility (2)

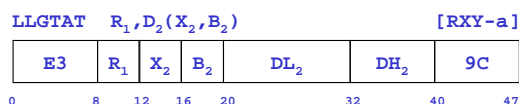
### LOAD HIGH AND TRAP



### LOAD LOGICAL AND TRAP



### LOAD LOGICAL THIRTY ONE BITS AND TRAP

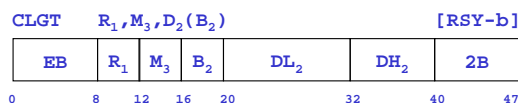
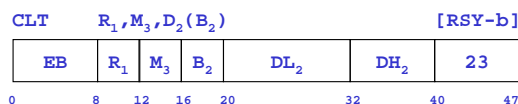


Continuing with the load-and-trap facility:

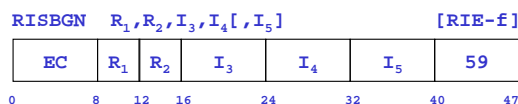
- LFHAT loads a 32-bit value into bits 0-31 of a register, and leaves bits 32-63 unchanged.
- LLGFAT loads a 32-bit value into bits 32-63 of a register, and bits 0-31 of the register are set to zero.
- LLFTAT loads a 31-bit value (bits 1-31 of the four-byte second-operand location) into bits 33-63 of a register, and bits 0-32 of the register are set to zero.

## Miscellaneous Instruction-Extensions Facility

- Provides storage-operand form of COMPARE LOGICAL AND TRAP:



- Provides non-condition-code-setting form of ROTATE THEN INSERT SELECTED BITS



- Facility indication bit 49

The miscellaneous instruction-extensions facility adds three instructions that are variations of instructions added in the System z10. This facility – as well as several others in the zEC12 – are indicated by facility indication 49 (as stored by STORE FACILITY LIST EXTENDED).

The z10 added the COMPARE LOGICAL AND TRAP instructions, each of which compared a value in a register with either another register or with an immediate field in the instruction. The two new forms of the instruction compare a value a register with a storage location; otherwise the operation is identical to the existing COMPARE LOGICAL AND TRAP instructions.

Both CLT and CLGT are of the RSY instruction format (subformat –b). The first operand contains either a 32-bit (CLT) or 64-bit (CLGT) value in a general register which is compared with a 4- or 8-byte operand in storage. The second operand is designated by a base register and a 20-bit signed displacement field. The trap conditions are specified by the 4-bit  $M_3$  field (similar to the branch mask of branching instructions). Note, HLASM provides extended mnemonics for these instructions, similar to the existing extended mnemonics for compare-and-trap instructions.

ROTATE THEN INSERT SELECTED BITS (RISBG) is one of the most powerful instructions in the CPU. Added in the z10, it provides the means of rotating and extracting bits from a register, and setting the condition code based on the result. Compiler- and system-development groups found that having the condition code set was not always optimal, so a separate instruction, RISBGN, was developed which does not set the CC. Otherwise, execution is identical to the original RISBG instruction.

## Transactional Execution – Overview

- **Statement of Problem being Addressed**
  - ▶ **Need for improved multiprocessing capabilities**
    - Issues with conventional MP serialization
  - ▶ **Speculative execution (e.g., for Java)**
    - Fall-back from partial in-lining that takes a side exit
    - Null checking that results in code scheduling barriers
- **Transactional execution**
  - ▶ **Controls**
  - ▶ **Instructions**
  - ▶ **Processing**
  - ▶ **Abort processing**
  - ▶ **Constrained transactions**

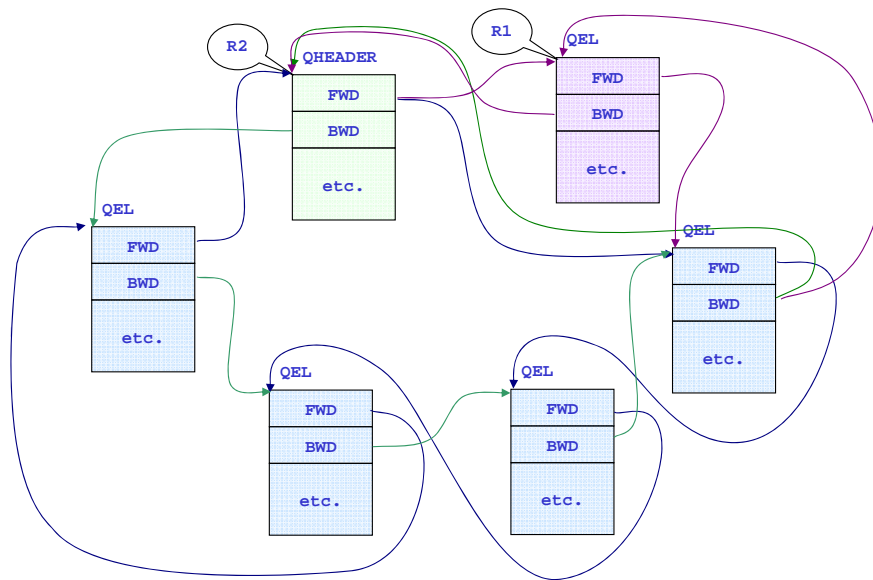
The majority of the slides in this presentation describe the transactional-execution facility (or TX facility, for short).

We begin with a description of two very different problems being addressed by the facility:

- The need for improved multiprocessing capabilities that address limitations in existing serialization, and
- The means by which a program can speculatively execute a code path, and – based on observed state of the program or exceptions encountered – efficiently withdraw said execution, making it appear as if the execution never occurred.

The following slides will discuss the controls, instruction, processing, abort handling, and a special form of TX called a constrained transaction.

## Example of a Serialized Operation: Element Insertion at the Head of a Doubly-Linked Queue



SHARE 120 – Session 12670

Use slide-show mode to make sense of this slide!

16

This slide illustrates a doubly-linked list into which a new queue element is to be inserted at the head of queue. The queue header (shown in green) and the existing queue elements (shown in blue) each contain a forward and backward pointer. Thus, in order to insert an element into the queue, multiple discontinuous storage locations must appear to be simultaneously updated (as observed by other CPUs and the channel subsystem) in order for the queue to retain integrity.

- The forward pointer of the queue header must be updated to point to the newly-inserted element.
- The backward pointer of original first element must be updated to point to the newly-inserted element.
- The forward pointer of the inserted element must point to original first element on the queue.
- The backward pointer of the inserted element must point queue header.

In order to maintain the integrity of the queue, the program will usually acquire some form of serialization such as a lock (also known as a semaphore or mutex). This technique is illustrated on the following slide.



## Example of a Serialized Operation: Sample Code Fragment using Locks

```

* R1 - address of the new queue element to be inserted.
* R2 - address of the insertion point (i.e., head of queue).

NEW    USING  QEL,1           Make new 1st QEL addressable.
HDR    USING  QEL,2           Make queue header addressable.
OLD    USING  QEL,3           Make old 1st QEL addressable.

      SETLOCK OBTAIN, ...     Serialize access to queue.
      LG      3,HDR.QEL_FWD   Point to original 1st element.
      STG     1,HDR.QEL_FWD   Update header's forward pointer.
      STG     1,OLD.QEL_BWD   Update orig. element's back ptr.
      STG     2,NEW.QEL_BWD   Update new element's backward ptr.
      STG     3,NEW.QEL_FWD   Update new element's forward ptr.
      SETLOCK RELEASE, ...

      ...

QEL    DSECT                   Common DSECT for header or QEL.
QEL_FWD DS    AD              Forward pointer.
QEL_BWD DS    AD              Backward pointer.
      DS     XL48              Queue element payload.

```

This assembler programming example shows how the update to these four storage locations can be accomplished using classic locking mechanisms.

Shown in the first highlighted section of code, the SETLOCK OBTAIN macro instruction is used to illustrate any number of locking, semaphore, or other serialization techniques that may be used to ensure that only one CPU is executing the following code fragment at any one time.

After obtaining the serialization, the program loads the address of the original first element on the queue into general register 3, and then proceeds to update the four key objects needed to insert the new element.

Finally, after performing the update, the SETLOCK RELEASE macro instruction (in the second highlighted section) illustrates the releasing of the serialization.

Note: In this example, because the forward and backward pointers appear in the same location in both the queue header and queue element, a single DSECT (QEL) is used.

## Problems with Conventional Serialization:

- **Coarse-grained locking**
  - ▶ Usually require serializing a much broader resource than what is being actually being accessed
  - ▶ With finer-grained serialization, multiple locks may be required
    - Hierarchy issues, potential dead-locks
    - Natural access may not lend itself to imposed hierarchy
- **Various recovery issues**
  - ▶ Lock cannot be acquired in a timely manner
  - ▶ Unexpected event encountered while locked

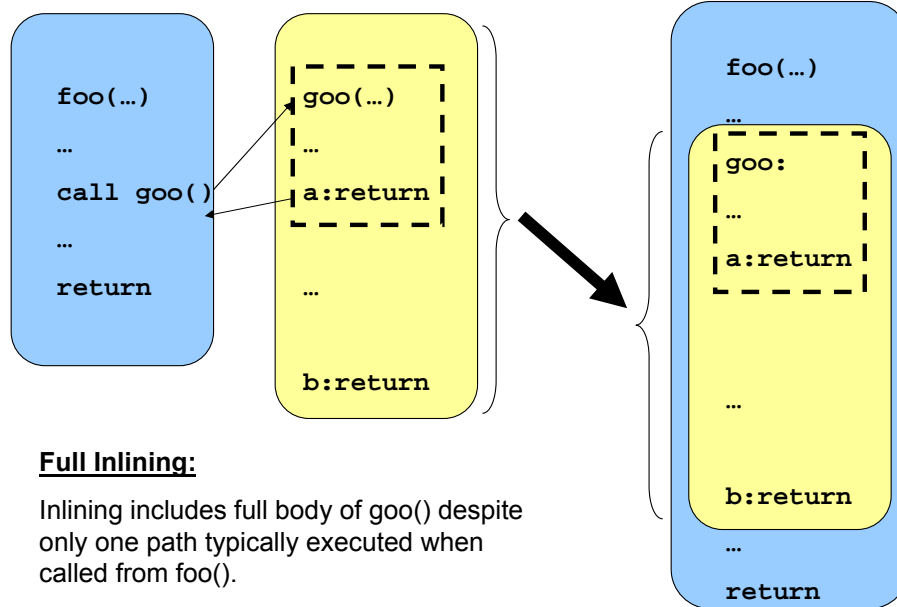
The problem with using classic locking techniques is that, in general, they serialize a much broader scope of resources than is actually being accessed. For example, in the queue illustration, the queue may contain millions of elements, yet only one is being updated. Even in a multiprocessing environment, many such data structures are serialized by coarse-grained lock, when it is rare to have multiple CPUs update the same location. (Nonetheless, the rare case must always be handled properly.)

Finer-grained serialization may exploit multiple levels of locking, but with such a hierarchy, there is the issue of potential deadlocks if the locks are not acquired and released in the proper sequence. Furthermore, finer-grained serialization imposes a regimen on the program that may be more complicated and error prone.

Additionally, the program must accommodate scenarios where either (a) the lock cannot be obtained, and (b) a task holding a lock encounters an unexpected condition (for example, abnormal end). Often the occurrence of “b” results in “a.”

Many years ago, IBM developed the PERFORM LOCKED OPERATION (PLO) instruction which provided a means by which separate storage locations could be updated under serialization provided by a configuration-wide lock. However, PLO does not co-exist well with classic forms of serialization that use compare-and-swap types of updates.

## Partial Inlining in Java

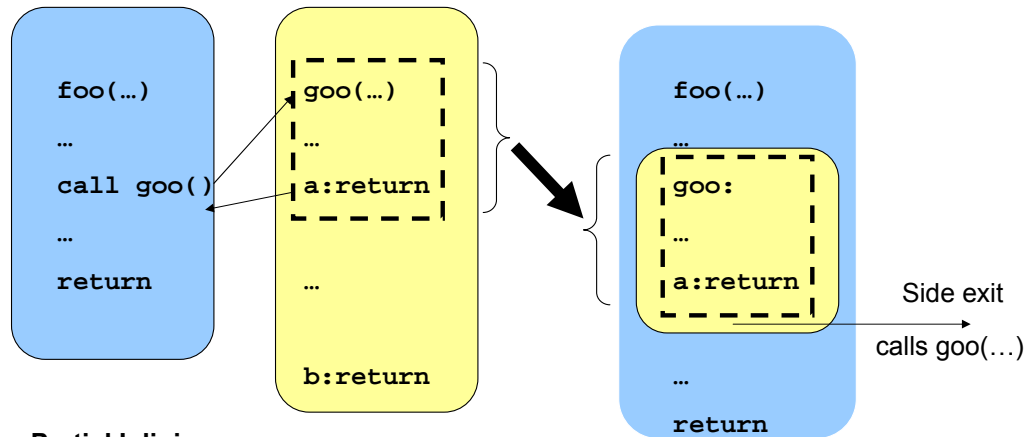


This slide illustrates a process that Java exploits to improve the performance of certain function calls.

On the left, we see the function foo() calling function goo(). Function goo() has a normal sequence of execution which returns at location "a," and an alternate sequence of execution which returns at location "b."

Java may restructure the functions such that a copy of the called function goo() is contained within the calling function foo(); this operation is called in-lining, and may be beneficial in minimizing instruction-cache references. If Java performs full in-lining of function goo(), it incorporates the function in its entirety inside the calling function foo(). This has the disadvantage of making foo() larger than it needs to be, increasing its instruction-cache footprint.

## Partial Inlining in Java



### Partial Inlining

Inlining includes only typically executed path in `goo(...)`.

Better i-cache footprint.

Fallback to calling `goo(...)` if alternative path taken.

In this slide, we illustrate the partial in-lining of function `goo()`. Only the commonly-executed sequence of code is placed into the in-lined version. This yields a better instruction-cache footprint.

However, if the execution of the partially-in-lined version of `goo()` determines that it needs to execute the alternate code sequence, it must call the full version of `goo()` to execute that code sequence.

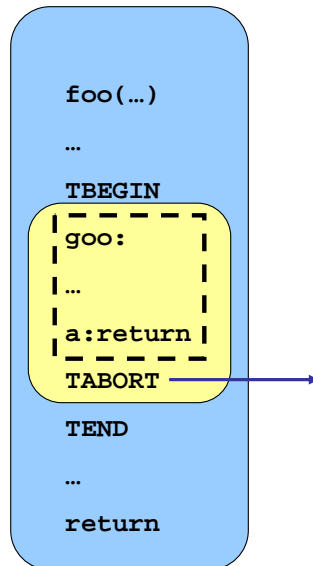
## Partial Inlining in Java



- Any changes to global state inside partial inlined path of `goo()` need to be undone if side exit taken
- Can be done in s/w, however overhead is prohibitive, opportunities are limited

The problem with partial in-lining of the function `goo()` is that if the alternate code path must be called, then any state changes made by the in-line version must be undone. This can be accomplished in software, however (a) the overhead is prohibitive, and (b) it significantly increases the complexity of code generation.

## Partial Inlining in Java



- Wrap partial region in TBEGIN/TEND
- TABORT rolls back state on side exit

This code fragment illustrates the bracketing of the partially-in-lined version of `goo()` with new instructions that are part of the transactional-execution facility: `TBEGIN` and `TEND`.

While executing the partially-in-lined version of `goo()` within a transaction, any changes to storage are not visible to other CPUs and the I/O subsystem until the `TEND` instruction completes.

Alternatively, if the function detects a situation which requires the calling of the full `goo()` function, it can execute a `TABORT` instruction to deliberately abort the transaction. In this case, all transactional stores made during transactional execution are withdrawn, as if they never occurred.

## General Speculation in Java

- **Java imposes implicit NULLCHKs on de-referenced pointers**
- **NULLCHKs are strongly ordered with respect to other global state changes and exception checks**
- **Strong ordering acts as a pipeline code scheduling barrier**

Another characteristic of Java is that it imposes strongly-ordered null checking on dereferenced pointers. That is, the checking must appear to occur before any use of the pointer is attempted.

This ordering does not necessarily lend itself to efficient scheduling of instructions in the pipeline. As we will see in the next slide, transactional execution provides a means of evading this strong ordering requirement.

## General Speculation in Java

```
do {
loop_start:
  o = o.next;
  codeA
  t = o.g;
  codeB
}
goto loop_start
...following code
```

JIT must insert implicit NULLCHK on O

→ acts as scheduling barrier

→ prevents codeA/codeB from scheduling together

NULLCHK is delayed to loop edge

→ codeA/B can schedule freely

→ if O is NULL, transaction will abort safely

```
TBEGIN
do {
loop_start:
  o = o.next;
  codeA
  t = o.g;
  codeB
}
if (o != null) goto loop_start
TABORT
...
All normal loop exits execute TEND
```

In the code sequence shown in the upper left, two code fragments, codeA and codeB are shown. In between the two fragments, a pointer “O” is dereferenced to fetch the member “g”.

Ordinarily, Java would have to insert a check of “O” to ensure it was not null, before attempting to execute codeB.

However, in the example on the lower right, Java can bracket this code sequence with a transaction. In this case, the checking for a null value of “O” can be deferred, such that codeA and codeB can be scheduled more efficiently. Subsequently, the value of “O” can be checked, and if null, cause the entire code sequence to be aborted.

The vast majority of the time, one expects that a non-null value of “O” will be used, thus the transaction will not be aborted. In the rare case of an abort, the abort-handler code can deal with the null pointer.



## Transactional-Execution (TX) Mode

- **New CPU state:**
  - ▶ Introduced in the zEnterprise EC-12 CPU architecture
  - ▶ Initiated by TRANSACTION BEGIN instruction
  - ▶ Ended by either:
    - Outermost TRANSACTION END (TEND) instruction
    - Transaction abort
- **While in the transactional-execution mode:**
  - ▶ All storage accesses by the CPU appear to be block-concurrent to other CPUs and the channel subsystem
  - ▶ Transactional store accesses are either:
    - Committed to storage when the outermost transaction ends normally (via TEND), or
    - Completely abandoned if the transaction is aborted

The transactional-execution (TX) facility adds a new CPU state to the processor – the transactional-execution mode.

TX mode is started by an outermost TRANSACTION BEGIN instruction. The term outermost is used, because transactional execution can be nested (as shown on the following slide). TX mode is ended by either of the following:

- An outermost TRANSACTION END instruction being executed, or
- The transaction being aborted.

While the CPU is in the TX mode, all storage accesses by the CPU appear to be block concurrent (that is, they happen all at once), as observed by other CPUs and by the I/O subsystem. These transactional stores are either (a) committed to storage and made visible to other CPUs when the outermost TRANSACTION END instruction completes, or (b) completely abandoned if the transaction is aborted.

## Nested Transactions

Nesting Depth	A	CSECT	B	CSECT	C	CSECT
0		...				
0		TBEGIN ...				
1		...				
1		BAS 14,B				
1				CSECT		
1				...		
1				TBEGIN ...		
2				...		
2				BAS 7,C		
2						CSECT
2						...
2						TBEGIN ...
3						...
3						TEND ...
2						...
2						BR 7
2				...		
2				TEND ...		
1				...		
1				BR 14		
1		...				
1		TEND ...				
0		...				

The diagram illustrates nested transactions across three control sections (A, B, and C). The nesting depth starts at 0. Section A begins a transaction (depth 1) with 'TBEGIN ...' and calls section B. Section B begins a transaction (depth 2) with 'TBEGIN ...' and calls section C. Section C begins a transaction (depth 3) with 'TBEGIN ...' and ends it with 'TEND ...'. Section B then ends its transaction with 'TEND ...'. Finally, section A ends its transaction with 'TEND ...'. Callouts indicate that the 'Outermost Transaction' is the one started by A, the 'Innermost Transaction' is the one started by C, and 'Transactional Stores Committed' occurs when A's transaction ends.

This slide illustrates the nesting of transactions, and the transaction nesting depth (TND).

Initially, when the CPU is not in the TX mode, the nesting depth is zero, as shown in the column on the left.

When control section A executes the TBEGIN instruction (that is, the outermost TBEGIN), the CPU enters the TX mode, and the transaction nesting depth (TND) is set to one. Control section A then calls control section B.

CSECT B also contains transactionally-executed code. When CSECT B executes its TBEGIN, the nesting depth is incremented to two. CSECT B then calls CSECT C.

As with B, CSECT C also contains transactionally-executed code. The execution of its TBEGIN instruction causes the TND to be incremented to three. CSECT C then ends its transaction with a TRANSACTION END (TEND) instruction. In this case, the CPU remains in the TX mode, but the nesting depth is decremented to two. CSECT C then returns to its caller.

CSECT B also executes a TEND instruction, causing the CPU to remain in the TX mode, but decrementing the nesting depth to one. CSECT B then returns to its caller.

CSECT A also executes a TEND instruction. This causes the nesting depth to decrement to zero, thus the CPU commits all stores made during transactional execution to memory and then leaves the TX mode.

## Transactional-Execution Controls

- **Control Register 0 (system controls):**
  - ▶ **Bit 8:** Transactional-execution control
    - Set by O/S to enable transactional execution
  - ▶ **Bit 9:** Program-interruption filtering override
- **Control Register 2 (task-related controls):**
  - ▶ **Bit 61:** Transaction diagnostic scope (prob / sup)
  - ▶ **Bits 62-63:** Transaction diagnostic control
- **TDBA:** Address of TX diagnostic block (TDB)
- **TAPSW:** PSW used if a transaction aborts
- **TND:** Transaction nesting depth
- **Facility indication bit 73**

There are several new controls affecting the transactional-execution facility:

Control register zero contains system-wide controls, as follows:

- Bit 8 indicates that the OS has enabled the TX facility. Because the facility requires OS support, this bit is set to zero by default, such that OS's that do not support TX can execute compatibly on a zEC12.
- Bit 9 is used in conjunction with program-interruption filtering, to be discussed later. The OS can override any program interruption filtering.

Control register two contains task-related controls used in the debugging of a transaction. These bits are set by various system-level debuggers.

- Bit 61 indicates the scope of control for the transaction-diagnostic control in bits 62-63 (that is, whether the control affects only problem state or both problem state and supervisor state).
- Bits 62-63 contain a diagnostic control. When nonzero, the diagnostic control causes various random aborts of transactions, thus allowing the testing of their abort handler.

An outermost TBEGIN instruction can specify the address of a transaction diagnostic block (TDB) into which various information is stored if the transaction is aborted. The address of this block is maintained in the TDB address (TDBA).

When an outermost transaction begins, it sets a transaction abort PSW (TAPSW) that is loaded if the transaction is aborted. More on this later.

The transaction nesting depth (TND) may be inspected by the program regardless of whether or not the CPU is in the TX mode.

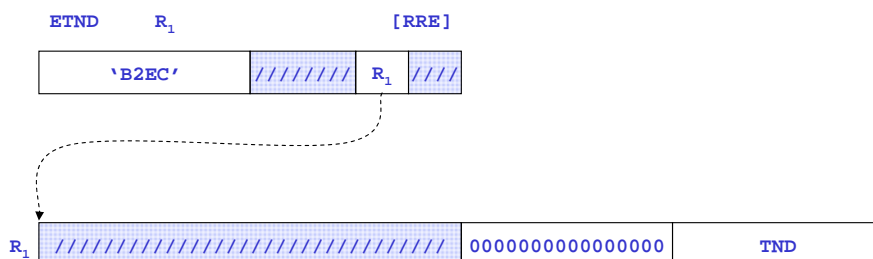
The availability of the TX facility is indicated by facility bit 73 (as stored by STORE FACILITY LIST EXTENDED). Note, even though a processor may provide the transactional-execution facility to logical partitions (LPARs), the facility may not be available when running in a virtual machine under the z/VM operating system.

## TX Facility Instructions

- **EXTRACT TRANSACTION NESTING DEPTH (ETND)**
  - ▶ Retrieves current nesting depth
- **NONTRANSACTIONAL STORE (NTSTG)**
  - ▶ Performs a store that will be committed regardless of whether the transaction aborts
  - ▶ Leaves “bread crumbs” for diagnostic purposes
- **TRANSACTION ABORT (TABORT)**
  - ▶ Deliberately causes a transaction to be aborted
  - ▶ User-defined abort code
- **TRANSACTION BEGIN**
  - ▶ Initiates (or continues) transactional execution
    - Nonconstrained transaction: `TBEGIN`
    - Constrained transaction: `TBEGINC`
- **TRANSACTION END (TEND)**
  - ▶ Outermost TEND causes transactional stores to be committed.

This slide summarizes the new instructions provided by the TX facility. Each of these will be described in detail in the following slides.

## EXTRACT TRANSACTION NESTING DEPTH



- Current transaction-nesting depth (TND) placed in bits 48-63 of register  $R_1$ 
  - Zero means not in transactional-execution mode
  - Bits 0-31 of  $R_1$  unchanged; bits 32-47 set to zeros
- Condition code is unchanged
- Exceptions:
  - PIC 0013 if transactional-execution control (CR0.8) is zero
  - PIC 0018 if issued in the constrained TX mode

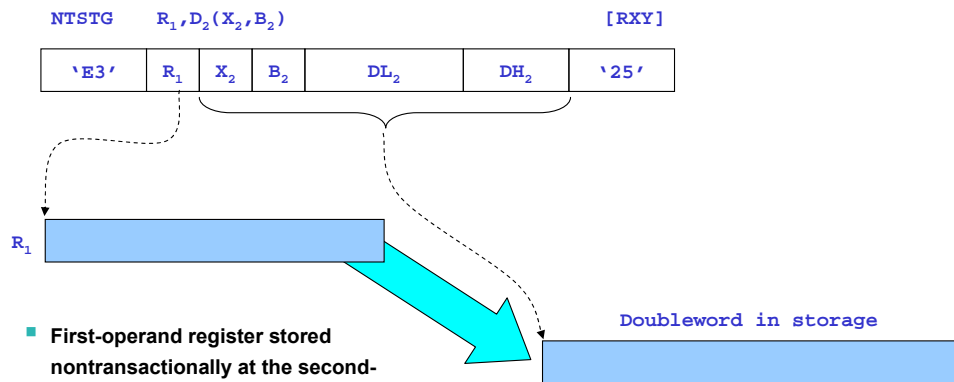
The EXTRACT TRANSACTION NESTING DEPTH (ETND) instruction provides a means by which the program can inspect the current transaction nesting depth (TND).

ETND is an RRE-format instruction with a single general register operand. The current nesting depth is placed in bits 48-63 of the register, and zeros are placed in bits 32-47; bits 0-31 of the register are unmodified. Also, the condition code remains unchanged.

As will be seen in most of the TX instructions, there is a special-operation exception recognized if the OS has not enabled the TX facility (in CR0.8).

Also described here is a new program exception, the transaction-constraint exception, if the instruction is executed while the CPU is in the constrained TX mode. More on nonconstrained versus constrained towards the end of this presentation.

## NONTRANSACTIONAL STORE



- First-operand register stored nontransactionally at the second-operand location
  - Leaves “bread crumbs” on abort
- Condition code is unchanged
- Exceptions:
  - PIC-0006 if 2<sup>nd</sup> operand not doubleword aligned
  - PIC-0018 if issued in constrained TX mode

The astute observer will have noticed that debugging transactional execution may prove to be challenging. This is because when a transaction is aborted, all evidence of transactional stores vanish.

As we will shortly see, there is the possibility of having some diagnostic information retained in registers, however additional information in storage would also be useful.

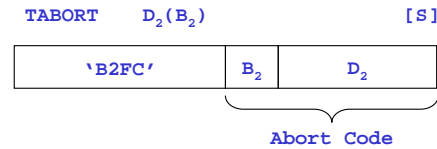
The NONTRANSACTIONAL STORE (NTSTG) instruction provides a means by which stores are retained following the abort of a transaction.

NTSTG is very similar to a regular STORE (STG) instruction. It is an RXY-format instruction, with the first operand designating a 64-bit register to be stored, and the second operand using a base, index, and long displacement to designate a storage location.

NTSTG differs from STG as follows:

- The storage operand must be on a doubleword boundary; otherwise, a specification exception is recognized.
- If NTSTG is executed in the constrained TX mode, a transaction-constraint exception is recognized. Regular STG can be used in the constrained TX mode.
- Stores made by NTSTG are retained if a transaction aborts; stores made by STG disappear.

## TRANSACTION ABORT



- Transactional execution aborted with code specified by second-operand address
- Condition code of the transaction-abort PSW (TAPSW) set as follows:
  - Bit 0 set to 1
  - Bit 1 set to bit 63 of the 2<sup>nd</sup>-operand address (abort code)
- Exceptions:
  - PIC 0003 if target of execute-type instruction
  - PIC 0006 if abort code < 256
  - PIC 0013 if transactional-execution control (CR0.8) is zero
  - PIC 0018 if issued in constrained TX mode

The TRANSACTION ABORT instruction provides the program with the means of deliberately aborting a transaction. It may be used, for example, to implement Java's partial in-lining or null-check reordering, described earlier.

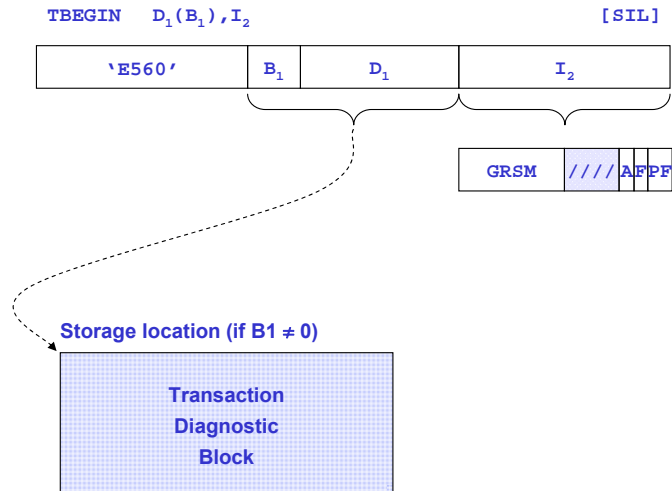
TABORT is an S-format instruction, having a second-operand address designated by a base register and 12-bit unsigned displacement. However the operand address is not used to access storage; rather, the address forms the transaction abort code that is stored into any transaction diagnostic block (TDB) designated by the outermost TBEGIN instruction.

TX architecture reserves abort codes 0-255 for use by the CPU. If the abort code in the second-operand address is less than 256, then a specification exception is recognized.

When a transaction is aborted, the condition code in the transaction-abort PSW indicates the likelihood of successful execution if the transaction is attempted again. CC2 means there is a potential for successful completion, and CC3 means that there is little potential for successful completion. TABORT sets bit 18 of the transaction-abort PSW to one, and bit 19 of the transaction-abort PSW is set to bit 63 of the second-operand address. Thus, the condition code is set to either 2 or 3, depending on whether the rightmost bit of the address is 0 or 1, respectively.

TABORT is unusual in that an execute exception (program interruption code 0003 hex) is recognized if TABORT is the target of an execute-type instruction (EX or EXRL).

## TRANSACTION BEGIN (TBEGIN)



The TRANSACTION BEGIN instruction is used to initiate or continue the execution of a *nonconstrained* transaction. We'll discuss nonconstrained versus constrained more in a later section.

TBEGIN is a SIL-format instruction having a storage operand designated by a base register and 12-bit unsigned displacement, and an immediate field containing various controls.

The operands are described further on the next slide.



## TBEGIN Operands

- When  $B_1 \neq 0$ , first operand of the outermost TBEGIN designates location of the transaction diagnostic block
  - ▶ Various diagnostic information stored on abort
- $I_2$  field contains various controls:
  - ▶ **GRSM**    **General-register (pair) save mask**
    - ◆ Designates GR pairs to be restored on abort
  - ▶ **A:**        **Allow AR modification control**
  - ▶ **F:**        **Allow floating-point-operation control**
  - ▶ **PIFC:**    **Program-interruption filtering control**
    - 0 – no filtering
    - 1 – limited filtering
    - 2 – moderate filtering

When an outermost TBEGIN is executed, and the base register of the first operand designates a register other than general register zero, the transaction-diagnostic-block address (TDBA) is set. If the transaction is aborted, diagnostic information will be (nontransactionally) saved in this block. If the  $B_1$  field of the instruction is zero, the TDBA is considered to be invalid, and no program-specified diagnostic information is saved.

The  $I_2$  field contains controls, as follows

The first 8 bits of the  $I_2$  field contain the general register save mask (GRSM). This field applies only to the outermost TBEGIN, and it is ignored for inner TBEGIN instructions. Each bit represents an even/odd pair of registers. If the bit is one, the contents of the register pair are preserved at the beginning of TX mode, and restored if the transaction is aborted. If the bit is zero, the register pair is neither saved nor restored.

General register are the only registers that may be saved, and if the transaction aborts, restored. Access registers and floating-point registers are not restored on an abort. The A and F controls, bits 12 and 13 of the  $I_2$  field, respectively, provide a means by which the program can prohibit a transaction from altering ARs or using FP instructions. Unlike the GRSM (which applies only to the outermost TBEGIN), the A and F controls assume effective values for each nested level of a transaction.

Finally, bits 14-15 of the  $I_2$  field contain a program-interruption filtering control. More on filtering in later slides.

## TBEGIN Processing

- **If TND = 0 (i.e., not in TX mode at beginning):**
  - ▶ If  $B_1 \neq 0$ , transaction-diagnostic-block address (TDBA) set from 1<sup>st</sup> operand address
  - ▶ Transaction-abort PSW set to next-sequential instruction address
  - ▶ General register pairs designated by  $I_2$  field are saved in model-dependent location
    - Not directly accessible by the program
- **Effective PIFC, A, & F controls computed**
  - ▶ Effective A = TBEGIN A & any outer A
  - ▶ Effective F = TBEGIN F & any outer F
  - ▶ Effective PIFC = max(TBEGIN PIFC, any outer PIFC)

This slide described the step-by-step processing that occurs during the execution of a TBEGIN instruction.

If the transaction-nesting depth (TND) is zero (that is, the CPU is not in the TX mode), then the following occurs:

1. If the  $B_1$  field designates a register other than zero, the transaction-diagnostic-block address (TDBA) is set. If the  $B_1$  field is zero, the TDBA is considered to be invalid.
2. The transaction-abort PSW is set to point to the instruction following the outermost TBEGIN.
3. Any general register pairs designated by the GRSM field of the instruction are saved in a model-dependent location that is not accessible by the program.

Regardless of whether the TND is zero, effective AR-modification, FR-modification, and program-interruption-filtering controls are computed, as follows.

- The effective A and F controls are the logical AND of the controls in the TBEGIN instruction and the respective controls in any outer TBEGIN instruction.
- The effective PIFC is the maximum of the control in the TBEGIN instruction and any outer TBEGIN instruction.

## TBEGIN Processing

- Transaction nesting depth (TND) incremented
- If TND transitions from 0 to 1, CPU enters the transactional-execution mode
  - ▶ Otherwise, CPU *remains* in transactional-execution mode
- Condition code set to zero
  - ▶ When instruction following TBEGIN receives control:
    - TBEGIN success indicated by CC0
    - Aborted transaction indicated by nonzero CC
- Exceptions:
  - ▶ Abort code 13 if nesting depth exceeded
  - ▶ PIC 0003 if target of execute-type instruction
  - ▶ PIC 0006 if either
    - PIFC is invalid (value of 3)
    - First-operand address not doubleword aligned
  - ▶ PIC 0013 if transactional-execution control (CR0.8) is zero
  - ▶ PIC 0018 if issued in constrained TX mode

Continuing with TBEGIN processing:

The transaction nesting depth (TND) is incremented. If the TND transitions from zero to one, the CPU enters the TX mode; otherwise, the CPU remains in the TX mode.

The condition code is set to zero.

Note, when a nonconstrained transaction aborts, control is passed to the transaction-abort PSW (TAPSW), the instruction address of which points past the outermost TBEGIN instruction. The condition code in the TAPSW will either be 2 or 3. Thus, the instruction following the outermost TBEGIN instruction is expected to be a conditional branch instruction. If the CC is zero, then it means that the transaction successfully initiated, and control is expected to fall through to the next instruction. If the CC is nonzero, then it means that the transaction was aborted, and control is passed to an abort handler.

As shown on this slide, there are various exception conditions. Of note:

- There is an abort condition if the maximum nesting depth of 16 is exceeded.
- A specification exception is recognized if a reserved PIFC value is coded, or if the first-operand address is not on a doubleword boundary.
- An execute exception is recognized if the instruction is the target of an execute-type instruction.
- If attempted in the constrained TX mode, a transaction-constraint exception is recognized.

## TRANSACTION END

TEND       $D_2(B_2)$       [S]



- If CPU is in the TX mode, transaction-nesting depth decremented
- If resulting TND is zero, CPU leaves the TX mode
  - All transactional and nontransactional stores are committed to storage
- Effective A, F, and PIFC controls returned to previous nesting-depth's values
- Condition code
  - 0 – CPU in TX mode at beginning of instruction
  - 1 –
  - 2 – CPU not in TX mode at beginning of instruction
  - 3 –
- Exceptions:
  - PIC 0003 if target of execute-type instruction
  - PIC 0013 if transactional-execution control (CR0.8) is zero

The TRANSACTION END (TEND) instruction is used to end a section of code that is executing in the TX mode. It is an S-format instruction, but has no operands.

If the CPU is in the TX mode, the transaction nesting depth (TND) is decremented.

If the resulting TND is zero, then the CPU leaves the TX mode, committing all stores to memory. Otherwise, the effective A, F, and PIFC control return to the value of the previous nesting-depth values.

TEND may be executed even if the CPU is not in the TX mode. The condition code indicates whether the CPU was in the TX mode (0) or not (2).

As with certain other TX-facility instructions, TEND may not be the target of an execute-type instruction.

## TX Facility Restricted Instructions

- All control and I/O instructions
- Any instruction that causes tracing or monitor-event counting
- SUPERVISOR CALL
- Any instruction that requires HW coprocessor assistance (e.g., message-security assists, compression call)
- Any instruction that is intercepted by the hypervisor
- Subject to the A & F controls (I<sub>2</sub> operand of TBEGIN):
  - ▶ Any instruction that modifies an access register
  - ▶ Any floating-point instruction.
  
- Transaction is aborted if a restricted instruction is executed
  - ▶ Abort code 11

The TX facility imposes numerous restrictions on the instructions that can be executed when the CPU is in the TX mode. Each of these are enumerated on this slide.

If a restricted instruction is attempted, the transaction is aborted with abort code 11 (see later slides for a complete list of abort codes).

## TX Operation

- **TX initiation**
  - ▶ Following successful execution of outermost TBEGIN, CPU enters the TX mode
    - Transaction-abort PSW set to address of next-sequential inst.
    - CC0 set to indicate successful execution.
- **TX normal ending**
  - ▶ As a result of successful execution of outermost TEND, all stores performed in TX mode are committed
    - Visible to other CPUs and channel subsystem
  - ▶ Registers are *not* restored!
- **TX abort**
  - ▶ Execution resumes at the TX-abort PSW
    - Nonconstrained TX:      Points past TBEGIN
    - Constrained TX:        Points at TBEGINC !!
  - ▶ General registers designated by outermost TBEGIN's GRSM are restored
    - Other GRs, all ARs and all FPRs are not restored!

This slide reviews the operation of the CPU while in the TX mode. Of particular note:

When a transaction completes normally, there is no restoration of general registers, access registers, or floating-point registers.

When a transaction aborts, registers designated by the GRSM field of the outermost TBEGIN instruction are restored; all other GRs, and all ARs and FRs and the floating-point control register (FPCR) are not restored.

## Reasons for a Transaction Abort

- **Any interruption**
  - ▶ External (2) ← text in parentheses is the abort code
  - ▶ Program (unfiltered, 4; filtered, 12)
  - ▶ Machine check (5)
  - ▶ I/O (6)
- **Conflict (with other CPU)**
  - ▶ Fetch (7)
  - ▶ Store (8)
- **Overflow**
  - ▶ Fetch (9)
  - ▶ Store (10)
- **Restricted instruction attempted (11)**
- **Nesting-depth exceeded (13)**
- **Cache-related**
  - ▶ Fetch (14)
  - ▶ Store (15)
  - ▶ Other (16)
- **Miscellaneous causes (255)**
  - ▶ Millicode timer, quiesce operation, SIE exit, operator intervention or SIGP equivalent, &c.)
- **TABORT instruction (> 255)**

This slide enumerates the various reasons that a transaction may be aborted. The numbers shown in the parentheses are the abort code.

Conflict conditions (codes 7, 8, and 14-16) represent cases where other CPUs have attempted to fetch from or store into locations that the transaction has stored into or fetched from, respectively.

Note that a transaction may be aborted due to any interruption, including an external interruption (such as a time-slice ending) or an I/O interruption. Thus, it is never assured that a transaction will complete on its first execution. Examples of re-driving an aborted transaction are shown below.

## TX Abort Action

- **Most transactional activity discarded:**
  - ▶ All transactional stores discarded
  - ▶ General registers designated by GRSM restored to values prior to outermost transaction
- **Stuff that persists following an abort:**
  - ▶ Nontransactional stores (NTSTG instruction[s])
  - ▶ ARs & FPRs retain any modifications
  - ▶ GRs not designated by GRSM retain any modifications
- **PSW set from TX-abort PSW**
  - ▶ Points to instruction following TBEGIN instruction
  - ▶ Condition code set to indicate reason
  - ▶ When aborted due to interruption, TX-abort PSW stored as interruption-old PSW (with CC indicating cause)
- **If outermost TBEGIN B<sub>1</sub> field is nonzero, TDB is stored**

When a transaction is aborted, all transactional stores are discarded ... not just as observed by other CPUs, but by the CPU executing the transaction as well. Thus, it's as if the stores never occurred. (There may be some lingering hints of stores having occurred, such as change bits remaining set in the storage keys, but as far as an application is concerned, there's no evidence.)

Additionally, any general register pairs that are specified in the general register save mask are restored to their pre-TX-mode contents. If all eight register pairs (that is, all 16 GPRs) are specified, then all contents of the registers are restored. Note, saving and restoring registers does consume CPU cycles, so specifying the minimum register set is recommended.

Things that are retained following a transaction's abort include any nontransactional stores made by the NTSTG instruction. Also, any ARs or FRs that were modified by the transaction retain their changes.

The transaction-abort PSW that was set by the outermost TBEGIN instruction, with a condition code set to indicate the severity of the abort, becomes effective. Thus, the instruction following the TBEGIN receives control following the abort of a nonconstrained transaction (or the TBEGINC instruction receives control following the abort of a constrained transaction ... more on this later).

Finally, if the B<sub>1</sub> field of the outermost TBEGIN instruction was nonzero, then a *program-specified* TDB is stored. Note, if the transaction is aborted due to program interruption conditions, then a program-interruption TDB is also saved in the prefix area of low storage.



## Abort Handling

- **For nonconstrained transaction, the instruction following the TBEGIN receives control**
  - ▶ **CC1 – should never occur (TDB became inaccessible after initial accessibility check)**
  - ▶ **CC2 – transient condition; re-execution of transaction may be successful**
    - **May want to limit number of re-drives**
  - ▶ **CC3 – persistent condition; re-execution of transaction not likely to be successful**
    - **May want to branch to fall-back path that uses conventional serialization techniques**
    - **NOTE: If fall-back path uses C&S-type locks, transaction should also check availability of lock!**

The condition code that is set in the transaction-abort PSW is always a nonzero value. Thus, the instruction following the outermost TBEGIN can determine whether it is being executed due to the successful execution of TBEGIN, or if it is the result of a transaction being aborted.

CC1 represents an extremely rare condition that should never occur. This indicates that the program-specified TDB (the accessibility of which is checked during the execution of the outermost TBEGIN) has become inaccessible during the execution of the transaction. This can only happen due to unexpected key changes made by the operating system – and should never occur.

CC2 represents a transient condition such as a conflict with another CPU. This condition is likely to be temporary, thus a repeated attempt at executing the transaction is likely to produce a successful completion.

CC3 represents a persistent condition, such as having exceeded the nesting depth or encountering an operation exception. Without program intervention, these conditions are not likely to go away on their own. Thus repeated attempts at executing the transaction are likely to continue to abort.

One special condition is worth mentioning. If TX is being used for lock elision (that is, to avoid using a lock), but the fall-back code path uses a lock word, then the transactional execution code path should test the accessibility of the lock to ensure that it's free. That way, the TX code path can co-exist with any non-TX code paths that attempt to serialize on the same resources.

## Transaction Diagnostic Block (TDB)

- **Program-Specified: When outermost TBEGIN B<sub>1</sub> field is nonzero, program-designated block in storage**
  - ▶ When outermost TBEGIN B<sub>1</sub> field is zero, no program-designated TDB is stored on abort
- **Program-exception conditions:**
  - ▶ For program interruptions, prefix TDB stored at real location 6,144 – 6,399 (1800-18FF hex)
  - ▶ For certain SIE interceptions due to program-interruption condition, interception TDB stored at location designated by bytes 488-495 of the state description

The transaction diagnostic block is a 256-byte area in memory. There are three types of TDBs, of which zero, one, or two may be stored when a transaction aborts.

- When the transaction diagnostic-block address (TDBA) is valid – that is, when the B<sub>1</sub> field of the outermost TBEGIN instruction is nonzero – then a program-designated TDB is stored if a transaction is aborted.
- For program interruptions, a TDB exists in the prefix area at locations 1800-18FF hex. A subset of the TDB information is saved in the prefix-area TDB (only a subset is saved, because some information in the TDB such as the program-interruption ID is already saved in other prefix-area locations).
- For conditions that cause the CPU to leave the interpretive-execution state (called SIE interceptions), an interception TDB is stored at a location designated by the hypervisor (that is, LPAR or z/VM).

## TDB Contents:

0	Format	Flags	Reserved	Trans Nest. Depth
8	Transaction Abort Code			
16	Conflict Token			
24	Aborted Transaction Instruction Address			
32	EAID	DXC	Reserved	Program Interruption ID
40	Translation Exception ID			
48	Breaking-Event Address			
56	Reserved			
112	Model-Dependent Diagnostic Information			
128	General Registers			
248				

The contents of the TDB are shown here.

- Format: When zero, the remaining fields of the TDB are unpredictable. When one, the fields are as shown in this slide.
- Flags: Bit 0 indicates that the conflict token (bytes 16-23) is valid. Bit 1 indicates that the CPU was in the constrained TX mode when the abort occurred.
- TND: The transaction-nesting depth at the time of the abort.
- Transaction Abort Code: The CPU-generated code or second operand of the TABORT instruction.
- Conflict Token: The logical address at which a conflict was detected. This field is valid only if flag bit 1 is one.
- Aborted-Transaction Instruction Address: The instruction address at which the abort was detected.
- EAID, DXC, program-interruption ID, translation-exception ID, and breaking-event address are only stored for transactions that are aborted due to program interruptions, and some of these fields are only stored for certain program interruptions. See the *z/Architecture Principles of Operation (SA22-7832-09)* for details on these fields.
- Model-Dependent Diagnostic Information: IBM internal-use diagnostic information.
- General Registers: The contents of all 16 general registers at the time of the abort.

## Transaction Diagnostic Controls

- **Task-specific controls in CR2**
  - ▶ **CR2.61 – Transaction-diagnostic scope (TDS):**
    - 0 – TDC applies regardless of whether CPU is in the problem or supervisor state
    - 1 – TDC applies only when CPU in the problem state
  - ▶ **CR2.62-63 – Transaction-Diagnostic Control (TDC):**
    - 0 – Normal operation; no random aborts
    - 1 – Abort every transaction at a random instruction, but before execution of outermost TEND
    - 2 – Abort random transaction at a random instruction
    - 3 – Reserved
- **Allows debugger to drive fall-back path**
  - ▶ **SCP interface TBD**

Control register 2 contains information that is unique to a task (sometimes called process or dispatchable unit). z/OS alters CRs for each task that is dispatched.

Two new fields are added to CR2 in support of debugging a transaction:

Bits 62-63 contain a transaction diagnostic control (TDC).

0 - means that transactional execution will not be randomly aborted ... at least, not due to the TDC.

1 - means that every transaction will be aborted at a random instruction, but before the TEND instruction.

2 - means that random transactions are aborted at random instructions.

The TDC allows a debugger to deliberately cause a transaction to be aborted, thus allowing the testing of the abort-handler fall-back code path (that is, the path branched to by a nonzero condition code in the instruction following the TBEGIN).

Bit 61 contains the transaction diagnostic scope (TDS). This bit controls the effectiveness of the TDC (in bits 62-63). When the TDS is zero, the TDC applies to both the supervisor and problem states; when the TDS is one, the TDC applies only to the problem state.

## Program-Interruption Filtering

- **Application program can request that certain program-exception conditions not result in interruption**
  - ▶ PIFC operand in the  $I_2$  field of TBEGIN
  - ▶ Effective PIFC is highest of all nested TBEGINs
    - 0 – no filtering; all exception conditions result in interruption
    - 1 – limited filtering
    - 2 – moderate filtering
- **Program can handle exception condition on abort ... without ESTAE, ESPIE, &c.**
  - ▶ Potentially useful in speculative execution
  - ▶ Caveat empor: Program may block interruptions necessary for its forward progress! (e.g., page-translation exception)
- **O/S can override TBEGIN-specified PIFC with PIFC-override control (CR0.9)**

One of the very powerful features of transactional execution is program-interruption filtering. By means of the program-interruption-filtering control (PIFC, bits 14-15 of the  $I_2$  field of the TBEGIN instruction), the program can request that certain classes of program-exception conditions that occur during TX not result in a program interruption. Rather, the transaction is simply aborted, and control is passed to the abort-handler transaction-abort PSW.

Thus, without establishing any elaborate recovery environment (such as using the z/OS ESPIE or ESTAE macros), the program can efficiently receive control if it causes an exception. This may be particularly useful in a speculative-execution environment, for example, the Java null-checking scenario. If an access exception is recognized dereferencing a pointer, then the transaction aborts, and the program can try a more conservative code path instead.

However, the program can cause program exceptions to be filtered that are otherwise necessary for it to make progress. For example, if page-translation exceptions are filtered, then the OS may not see the exception and cause the page frame to be migrated in from auxiliary storage. Therefore, a transaction's fall-back path may need to reference storage locations that cause protection or translation exceptions in order to allow the OS to resolve the exceptions.

The OS has the ability to completely override program-interruption filtering by means of the program-interruption filtering override – bit 9 of control register 0.

## Example of a Transaction used for Lock Elision: Sample Code Fragment

\* R1 - address of the new queue element to be inserted.  
\* R2 - address of the insertion point; new element is inserted before the element pointed to by R2.

```
NEW    USING QEL,1      Make new 1st QEL addressable.
HDR    USING QEL,2      Make queue header addressable.
OLD    USING QEL,3      Make old 1st QEL addressable.
```

LHI	R15,10	Load retry count.
LOOP	TBEGIN TDB,X'C000'	Begin transaction (save GRs 0-3)
JNZ	ABORTED	Nonzero CC means aborted.
LG	3,HDR.QEL_FWD	Point to original 1st element.
STG	1,HDR.QEL_FWD	Update header's forward pointer.
STG	1,OLD.QEL_BWD	Update orig. element's back ptr.
STG	2,NEW.QEL_BWD	Update new element's backward ptr.
STG	3,NEW.QEL_FWD	Update new element's forward ptr.
TEND		End transaction.

...

ABORTED	JC	B'0101',NO_RETRY	CC1 or CC3 - don't bother retrying.
	JCT	R15,LOOP	Retry transaction a few times.
	J	NO_RETRY	No joy after 10x; do it the hard way.

This assembler program fragment illustrates the same code sequence as originally shown on slide 17.

In this example, the SETLOCK macro instructions are replaced with TBEGIN and TEND instructions, as shown in the first two highlighted areas. The code fragment uses general register 15 as an abort counter.

If the transaction is aborted, the code branches to the label ABORTED, where a determination is made as to whether the transaction should be re-driven. For condition codes 1 or 3, there is little chance of recovery, so the code is not reattempted (it branches to NO\_RETRY). However, if the condition code is 2, the count in general register 15 is decremented, and, if nonzero, the transaction is attempted again.

## General Comments

- **If the transaction is used for lock elision, and the fall-back path uses a lock, the transaction must (at least) fetch the lock word to see that it's available.**
  - ▶ Ensures that the transaction aborts if another CPU accesses the lock non-transactionally.
- **Coding of both transactional and fall-back path adds to the complexity of the code**
  - ▶ Hence, constrained transactions for small updates (see next slides)
- **If transactions are nested, outermost transaction must account for unanticipated abort conditions that occur in the inner transaction**
  - ▶ E.g., modified, but unrestored GRs, ARs, FPRs

As noted earlier, a transaction that is used for lock elision – but has a fall-back path or other code paths that use a conventional lock – should check to see if the lock is available, and if not, end and branch to an abort handler. This ensures that the transactional and conventional code can successfully co-exist.

Also, coding of a transaction may be quite simple, but having both transactional and fall-back code paths may be more complicated and require additional testing. Constrained transactions, to be discussed shortly, address this limitation by eliminating the fall-back path completely.

When transactions are nested, an outermost transaction might encounter aborts due to unanticipated conditions discovered in an inner transaction. For example, an outermost transaction might set the GRSM to save registers 0-3, saving only those registers that it modified. However, if this code called some other function – perhaps a library routine – that modified other registers, and then the transaction aborted, the program may be unprepared to deal with other changed registers that were altered by the called program.

## Constrained Transaction

- TX mode without fall-back path
  - ▶ Intended for extremely compact function (à la PLO)
  - ▶ No fall-back path!
  - ▶ Constrained TX facility indicated by facility bit 50
- Constraints:
  - ▶ No more than 32 executed instructions
  - ▶ No more than 256 bytes of instruction text
  - ▶ No nesting! Maximum transaction-nesting-depth = 1
  - ▶ No looping!! Forward relative-branching instructions only
  - ▶ Many more restricted instructions
    - All instructions implemented in Millicode
    - All SS-format instructions (i.e., no MVC, CLC, &c.)
    - All floating-point instructions
  - ▶ No more than four octowords of storage accessed (4x32)
  - ▶ No operand accesses to instruction stream
  - ▶ No address-space remapping
  - ▶ Integral-boundary alignment requirements

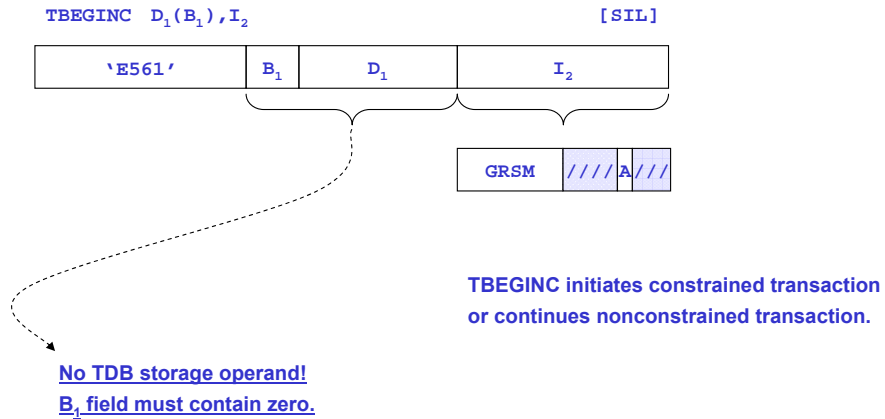
Coding a fall-back code path introduces a fair amount of complexity into transactional execution. The constrained transaction minimizes that complexity by eliminating the need for a fall-back path. However, the constrained transaction has significant additional restrictions (constraints), as enumerated on this slide.

Even though a constrained transaction may initially abort, it is assured of eventual completion.

Thus far, when we have discussed transactional execution, we have been discussing the *nonconstrained* TX mode (as initiated by a TBEGIN instruction). The following discussion describes an additional *constrained* TX mode (as initiated by the TBEGINC instruction).



## Constrained Transactions TRANSACTION BEGIN (TBEGINC)



This slide illustrates the TBEGINC instruction, a variant of TRANSACTION BEGIN that is used to initiate the constrained transaction. TBEGINC is similar to TBEGIN, except as follows:

1. There is no abort handler!
2. Since there is no abort handler, there is no need for a program-specified transaction-diagnostic block (TDB). The  $B_1$  field of the TBEGINC instruction must contain zero!
3. The controls in the  $I_2$  field are limited: there are no floating-point control (F) or program-interruption-filtering control (PIFC) bits.

## Constrained Transactions TRANSACTION BEGIN (TBEGINC)

- **If already in the constrained-TX mode, transaction-constraint exception recognized**
- **If current TND > 0, execution proceeds as if nonconstrained transaction**
  - ▶ **Effective F control set to zero**
  - ▶ **Effective PIFC is unchanged**
  - ▶ **Allows outer nonconstrained TX to call service function that may or may not use constrained TX.**
- **If current TND = 0:**
  - ▶ **Transaction-diagnostic-block address is invalid**
    - No instruction-specified TDB stored on abort
  - ▶ **Transaction-abort PSW set to address of TBEGINC!**
    - Not the next sequential instruction
  - ▶ **General-register pairs designated by GRSM saved in a model-dependent location not accessible by program**

Execution of TBEGINC is as follows.

1. Nesting within a constrained transaction is not permitted. If the CPU is already in the constrained TX mode, then a transaction-constraint program interruption is recognized (program interruption code 0018 hex).
2. If the transaction nesting depth is already greater than zero (meaning the CPU is in the nonconstrained TX mode), then execution simply proceeds as if this was a nonconstrained transaction. In this case, the effective F control is zeroed, and the effective PIFC remains unchanged. This allows an outer, nonconstrained transaction to call a service function that may or may not use constrained TX mode.
3. If the current TND is zero, then:
  1. The TDBA is marked as invalid (since there is no abort handler, there is no need for a transaction diagnostic block)
  2. The transaction-abort PSW is set to point directly at the TBEGINC instruction! This means that if the transaction is aborted, it will be re-driven – without attempting to branch to an abort handler.
  3. Any general registers specified by the GRSM are saved.

## Constrained Transactions TRANSACTION BEGIN (TBEGINC)

- **Effective A = TBEGINC A & any outer A**
- **TND incremented**
  - ▶ If TND transitions from 0 to 1, CPU enters the constrained TX mode
  - ▶ Otherwise, CPU remains in the nonconstrained TX mode
- **Instruction completes with CC0**
- **Exceptions:**
  - ▶ Abort code 13 if nesting depth exceeded
  - ▶ PIC 0006 is B<sub>1</sub> field is nonzero
  - ▶ PIC 0013 if transactional-execution control (CR0.8) is zero
  - ▶ PIC 0018 if issued in constrained TX mode

Continuing with TBEGINC execution:

1. The effective A control (allowing AR modification) is set to the control on the TBEGINC instruction logically ANDed with any outer value of A.
2. The transaction nesting depth is incremented. If the result is 1, then the CPU enters the constrained TX mode; otherwise, the CPU remains in the nonconstrained TX mode.
3. The instruction completes with CC0.

Various exception conditions may be recognized, as listed on this slide. Note, although a constrained transaction may only have a nesting depth of 1, the instruction recognizes an abort if the maximum TND is exceeded. This can occur if the CPU is already in the nonconstrained mode when TBEGINC is executed.

## Constrained Transaction

- **Abort conditions in constrained transaction:**
  - ▶ **Abort PSW points to TBEGINC instruction!**
    - NOT the instruction following it!
    - Abort condition causes entire TX to be re-driven!
      - No fail path!!
    - Re-execution of TBEGINC resets the condition code, thus branch immediately following TBEGINC is not productive
  - ▶ **CPU takes special measures to ensure successful completion on re-drive**
  - ▶ **Assuming no persistent conflict, interrupt, or constraint violation, the transaction is assured of eventual completion.**
- **Constraint violation:**
  - ▶ **PIC 0018 – indicates violation of transaction constraint**
  - ▶ **Or, transaction runs as if non-constrained**
  - ▶ **ASMAXCTX (HLASM exit) warns of static constraint violations**

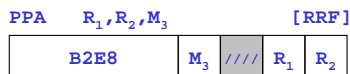
The astute reader may ponder, what prevents my CPU from looping forever in a constrained TX ... aborting and re-driving forever?

If any of the constraints are violated, this is possible. So, the constrained transaction must be carefully coded. HLASM provides an exit, ASMAXCTX, that helps in identifying constraint violations.

However, in the absence of constraint violations, there may still be causes for an abort, for example, if multiple CPUs repeatedly access the same storage locations. The CPU takes special measures to ensure that the re-drive of an aborted constrained transaction will eventually complete.

## Processor-Assist Facility (1)

### ■ Interface for program to request CPU assistance



- ▶ Facility indication bit 49
- ▶  $M_3$  field indicates which processor assist to perform
  - Currently, only the transaction-abort assist implemented ( $M_3 = 1$ )
    - May be invoked after *nonconstrained* transaction is aborted.
    - Bits 32-63 of general register  $R_1$  contain abort count.
    - Bits 0-31 of general register  $R_1$  and all of general register  $R_2$  ignored.
    - CPU takes actions to improve likelihood of successful re-execution of nonconstrained transaction
- ▶ Processor assist otherwise acts as a no-op
  - No visible change to the conceptual sequence

What about nonconstrained transactions? Is there any CPU tricks that can help them to complete successfully?

The processor-assist facility provides an instruction that can be used following the abort of a nonconstrained transaction. It provides a hint to the CPU as to how many times the transaction has aborted, and the CPU may take certain measures to help ensure that a subsequent execution of the transaction will succeed.

PPA is an RRF-format instruction designed to accommodate future assist functions; initially, only the transaction-abort assist is implemented, as indicated by an  $M_3$  value of 1.

Bits 32-63 of the general register designated by the  $R_1$  field contain a program-specified count of how many times a nonconstrained transaction has been aborted; initially, bits 32-63 of this register should contain zero. Bits 0-31 of general register  $R_1$  and all of general register  $R_2$  are ignored.

Other than the actions described above, PPA acts as a no-operation.

Note, the program should not attempt to “cheat” the recovery action by specifying a higher abort count than has actually occurred. A higher count does not necessarily provide a higher chance of successful redrive, and it may hurt program performance.

## Processor-Assist Facility (2)

### ■ Example of use:

```

LHI    15,0           Zero counter.
LOOP   TBEGIN TDB,X'F000' Restore GRs 0-7 if aborted.
       JNZ    ABORTED   Branch if aborted.
       :
       :               Transactional-execution code
       :
       TEND           End of transaction

ABORTED JC    B'0101',NO_RETRY CC 1 or 3; not worth retrying.
        AHI    15,1       Increment counter.
        CIJNL  15,6,NO_RETRY Give up after 6 attempts.
        PPA    15,0,1     Request CPU assistance.
        J     LOOP       Once more ... with feeling.

NO_RETRY DS    0H

```

This slide shows an assembler program's use of the PPA instruction. The first block of highlighted code shows the transactional-execution code; prior to this block, general register 15 is zeroed.

In the ABORTED section of code, we again check the condition code to see if there is a chance of recovery; if not, the code branches to NO\_RETRY. Next, the count in general register 15 is incremented, and if it has exceeded a threshold (six, in this case), we also branch to NO\_RETRY. Otherwise, the PPA instruction is executed, as shown in the second section of highlighted code, and the program branches back to attempt to re-execute the transaction.

## Transactional-Execution Summary

- **Benefits:**
  - ▶ **Efficient means of updating multiple, discontinuous objects in memory**
    - Without classic (coarse-grained) serialization such as locking
    - Potential for significant MP performance improvement
  - ▶ **Speculative execution without onerous recovery setup**
  - ▶ **Constrained TX for simple, small-footprint updates**
    - As long as constraints met, assured of eventual completion
- **Caveats:**
  - ▶ **Nonconstrained transaction may require more development / testing**
    - Fall-back path required to accommodate aborted TX
  - ▶ **Constrained TX has significant limitations**
    - If coded incorrectly, code could loop forever

Transactional execution has enormous potential in reducing overhead associated with classic coarse-grained serialization such as locking. Already IBM lab tests have shown extremely promising improvement in certain workloads that would otherwise suffer MP effects with more processors.

Transactional execution also holds promise in implementing speculative forms of execution, and may simplify many coding paths that would otherwise need cumbersome recovery environments (like ESTAE and ESPIE).

For quick multiple-location updates with a small memory footprint, the constrained transaction provides an effective means of high-performance processing.

However, there may be additional development required to develop and test unconstrained transactions, and constrained transactions have very stringent coding requirements.

## Enhanced-DAT Facility 2

- **System z10 introduced the enhanced-DAT facility**
  - ▶ **Format-1 segment-table entry (STE) defines 1 M-byte absolute segment frame**
  - ▶ **Access-control and fetch-protection bits in STE**
  - ▶ **Change-bit override in STE**
  - ▶ **Extends DAT-protection controls into region-table entries**
  - ▶ **Now called enhanced-DAT facility 1 (EDAT-1)**
- **Enhanced-DAT facility 2:**
  - ▶ **Extends EDAT-1 concepts to region-third table (RTTE)**
    - **Format-1 RTTE defines 2 G-byte absolute region frame**
    - **Access-control and fetch-protection bits in RTTE**
    - **Change-bit override in RTTE**
    - **Common-region control in RTTE**
  - ▶ **Facility indication bit 78**
  - ▶ **TLB redefined as a fully hierarchical structure**

The System z10 introduced the enhanced-DAT (EDAT) facility, now called the enhanced-DAT facility 1 (or just EDAT-1, for short). The features introduced by this facility are enumerated on this slide, however the one of most attention is the 1 megabyte segment-frames (often – but inaccurately – called large pages).

Enhanced-DAT facility 2 builds upon EDAT-1, providing a super-large 2 G-byte region frame. Similar to the changes made by EDAT-1 to the segment-table entry, EDAT-2 adds new controls to the region-third table entry, providing the format control, access- and fetch-protection controls and corresponding validity indication, change-bit override, and common-region controls.

Facility indication bit 78 designates the presence of the EDAT facility 2.

Additionally, with the introduction of this facility, the of the translation lookaside buffer (TLB) is redefined to be a completely hierarchical structure. While retaining compatibility with the former TLB structure consisting of page-table entries (PTEs) and common-region-and-segment-table entries (CRSTEs), the new structure allows for more flexible future design.



## Enhanced-DAT Facility 2: Region-Third-Table Entry (RTTE)

### ■ Format-0 RTTE:



### ■ Format-1 RTTE



**Explanation:**

ACC	Access-control bits	I	Invalid bit
AV	ACC/F validity bit	P	DAT-protection bit
CO	Change-bit override	TF	Table offset
CR	Common-region	TL	Table length
F	Fetch-protection bit	TT	Table type
FC	Format control		

This slide illustrates the changes to the region-third-table entry (RTTE) with EDAT-2.

The RTTE is extended to include a format control in bit 53. When bit 53 is zero, the definition of the RTTE is as originally defined in z/Architecture, shown in the upper illustration (format-0 RTTE).

When the FC is one, the definition of the RTTE is as shown in the lower illustration which includes the following:

- Bits 0-32           - Region-frame absolute address
- Bit 47             - ACC and F bit validity indication
- Bits 48-51        - Access-control bits for the region
- Bit 52             - Fetch-protection bit for the region
- Bit 53             - Format control
- Bit 54             - Region-protected bit
- Bit 55             - Change-bit override
- Bit 58             - Invalid bit
- Bit 59             - Common-region bit
- Bits 60-61        - Table type (01 binary)

## Enhanced-DAT Facility 2: CRDTE

### ■ COMPARE AND REPLACE DAT TABLE ENTRY

CRDTE  $R_1, R_3, R_2[, M_4]$  [RRF]

B98F	$R_3$	$M_4$	$R_1$	$R_2$
------	-------	-------	-------	-------

- ▶ **Compares doubleword designated by second operand with the contents of general register  $R_1$  (compare value)**
  - If equal:
    - Designated doubleword replaced by the contents of  $R_1+1$  register (replacement value)
    - TLB purged of at least all entries matching replaced value.
    - When the  $R_3$  field designates a register other than zero, TLB clearing limited to ASCE designated in the register
    - $M_4$  field contains local-TLB-clearing control
      - ◆ When zero, all TLBs on all CPUs cleared
      - ◆ When one, only the TLB of the CPU executing the instruction is cleared
  - If not equal, the doubleword designated by the second operand is loaded into general register  $R_1$
- ▶ **Improved performance when replacing live DAT-table entry**

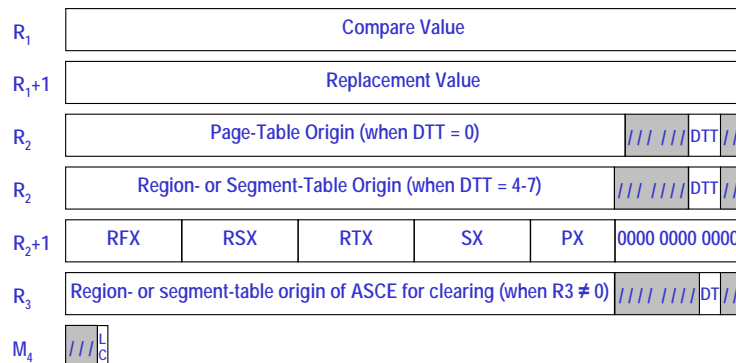
EDAT facility 2 also introduces the COMPARE AND REPLACE DAT TABLE ENTRY (CRDTE) instruction. CRDTE combines the storage-update operations of a compare-and-swap operation with the TLB purging that is associated with INVALIDATE DAT TABLE ENTRY or INVALIDATE PAGE TABLE ENTRY.

Normally, when an operating system needs to update the DAT table entry of an attached DAT table (that is, a table that may actively be used by the CPU for translation), then it must first invalidate the table entry and purge any TLB entries formed from it, make the necessary changes, and revalidate the entry. This assures that stale entries in the TLB won't accidentally be used by other CPUs, while one CPU makes the changes.

CRDTE combines these functions into a single instruction, improving the performance of such updates for all levels of DAT-table entries.

## Enhanced-DAT Facility 3: CRDTE

■ **CRDTE operands (similar to IDTE)**



**Explanation:**

DT	Designation type	RSX	Region 2 <sup>nd</sup> index
DTT	Designated table type	RTX	Region 3 <sup>rd</sup> index
LC	Local-clearing control	SX	Segment index
RFX	Region 1 <sup>st</sup> index	PX	Page index

The operands to CRDTE are similar to those used by COMPARE AND SWAP, INVALIDATE DAT TABLE ENTRY, and INVALIDATE PAGE TABLE ENTRY, as shown in this slide.

The first operand comprises an even / odd general register pair containing the compare value and the replacement value for a DAT-table entry designated by the second operand.

The second operand contains an even / odd general register pair designating the location of the DAT-table entry to be replaced. The even-numbered register contains the base address of the table, and a designated-table-type (DTT) indication. The odd-numbered register is in the form of a virtual address; based on the DTT, the appropriate portion of the virtual address (RFX, RSX, RTX, SX, or PX) is used to locate the entry in the table.

Assuming the comparison is equal, the table entry is replaced using a block-concurrent interlocked update, and the TLBs of all CPUs are cleared of at least that entry and any subordinate entries.

When the  $R_3$  field is nonzero, the clearing can be restricted to a particular ASCE

Additionally, CRDTE contains a  $M_4$  field containing a local-clearing control. When this bit is one, clearing is restricted to the CPU on which the instruction executes. This may avoid disruptive clearing in a uniprocessor environment or in an MP where an address space is only dispatched on a single CPU. More on local clearing in the following slides.

## Local-TLB-Clearing Facility

- Adds local-clearing control to TLB-clearing instructions:

- ▶ INVALIDATE DAT TABLE ENTRY (IDTE)

IDTE  $R_1, R_3, R_2[, M_4]$  [RRF]

B98E	R <sub>3</sub>	M <sub>4</sub>	R <sub>1</sub>	R <sub>2</sub>
------	----------------	----------------	----------------	----------------

- ▶ INVALIDATE PAGE TABLE ENTRY (IPTE)

IPTE  $R_1, R_2[, R_3[, M_4]]$  [RRF]

B221	R <sub>3</sub>	M <sub>4</sub>	R <sub>1</sub>	R <sub>2</sub>
------	----------------	----------------	----------------	----------------

- ▶ Local-clearing control in bit 3 of the *optional* M<sub>4</sub> field:

0 – perform global clearing (i.e., business as usual)

1 – perform local clearing (i.e., only on the CPU executing the instruction)

- ▶ Facility indication bit 51

– Local-clearing control ignored if facility not installed

The INVALIDATE DAT TABLE ENTRY and INVALIDATE PAGE TABLE ENTRY instructions have both been enhanced to provide a local-clearing control.

In the respective instruction images shown above, this is the M<sub>4</sub> field of the instruction. In the assembler syntax, this field is optional, and is effectively zero if the field is not coded.

The local-clearing control operates as described on the previous slide for COMPARE AND REPLACE DAT TABLE ENTRY: if the bit is zero, global clearing occurs on all CPUs; if the bit is one, local clearing on the CPU executing the instruction occurs.

If the local-TLB-clearing facility is not installed (as indicated by facility indication bit 51 being zero), then setting the LC bit to one on either of these instructions is ignored.

## zEnterprise EC12 CPU-Facility Summary

- **Improved facilities for multiprocessing**
  - ▶ Interlocked-access facility 2
  - ▶ Transactional-execution facility
  - ▶ Processor-assist facility
- **Improvements for programs using zoned & packed data**
  - ▶ DFP zoned-conversion facility
- **Various CPU facilities for code optimization**
  - ▶ Execution-hint facility
  - ▶ Load-and-trap facility
  - ▶ Miscellaneous general instructions
- **Improved DAT function**
  - ▶ Enhanced-DAT facility 2
  - ▶ Local-TLB-clearing facility
- **Potential for substantial performance improvements in selected workloads**

The IBM zEnterprise EC12 system introduces several powerful new CPU facilities as enumerated on this slide. The most significant of these changes is the transactional-execution facility which provides a game-changing paradigm in multiprocessor serialization, and has the potential of providing significant performance improvement for selected workloads. The interlocked-access facility 2 makes it significantly easier to manipulate shared data – without the bother of using locks or compare-and-swap type of operations.

Other enhancements provide improvements in legacy workload handling (DFP zoned-conversion facility), pipeline optimization (execution hint facility and miscellaneous-general-instructions facility), data validation (load-and-trap facility), and virtual-storage management (EDAT-2).

In combination, these facilities set the stage for improved CPU throughput.



For those in the audience, (a) give yourself a pat on the back for enduring an intense hour of discussion on the zEC12, and (b) I will gladly entertain any questions now.

For those reading this from the SHARE web site, if you have further questions, you are welcome to e-mail me at [dgreiner@us.ibm.com](mailto:dgreiner@us.ibm.com).