

Vol 5



**College Ruled White Paper**  
**Single Subject**

Dennison National Company, Holyoke, MA 01041

**31-987**

**80 Sheets/11 x 8½**

## VM/SP Installation

11-17-85

Claudia Dewey

EastFSC - route prt destination

VM Overview

→ used to be standalone

1966 CMS (Cambridge) and CP/40 were merged.

1967 Released w/o support, (Rel. 6. - public domain)

1972 VM/370 - licensed

1980 VM/SP

1982 HPO

VM under VM cannot test AP/MP code.

Components:

CP - DMK prefix

CMS - DMS

IPCS - DMM (IPCS) or DTV (IPCS/E)

With ver. 4 only IPCS/E is available, and  
is no longer an extra cost item.

RSCS - DMT

Distribution

Starter System - 3 tapes:

1. Object

2. PUT

3. Actual system (default addresses)

Pre-generated SIPO or SIPO/E (custom).

Two-year backlog.

Merged Product Tape

## Architecture

Jim Turner

CP runs disabled with translate off, in way 0.  
Only the HPO module which moves pages below the line  
can run with translate on.

CMS runs in a 360 environment using a BC mode PSW,  
(360 machines only had 5 channels).  
PSW actually held masks for each channel!  
Also, the actual interrupt code!

> TOD clock synchronization - a continuous duty?  
TOD - 64 bits  
bit 51 incremented each microsecond  
bit 31 changes each 1.0486 second  
The TOD clock will roll after 143 years.

CPU Timer decrements only while CPU is running.  
used for timing.

Interval Timer - decrements only while CPU is running.  
used to time dispatch time slice end.

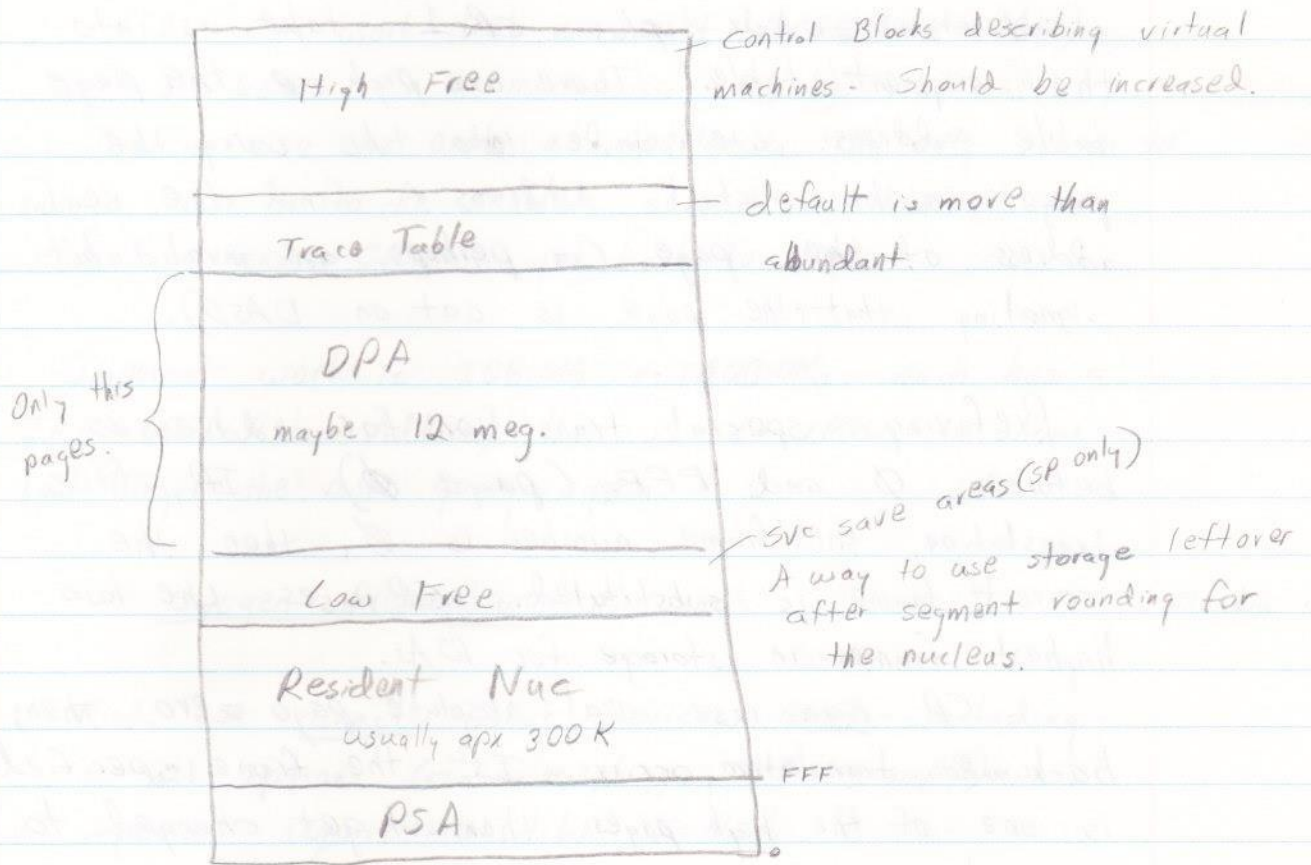
Restart interrupt causes CP to dump.

> Foil ARCD120 - what would prevent an I/O  
interrupt or machine check from occurring  
with PSW bit 14 on?

VMCF + IUCV signal via External Interrupt.

VM/SP uses SIOF + deferred condition codes.

DAT takes an average of 64 machine cycles. TLB cuts this dramatically.



Each segment table entry maps 64K. Last 3 bits are 0 + so may be used as flags.

29 = seg. protect (HPO)

30 = common segment

31 = segment invalid

Page table has 16 entries ( $16 \times 4K$ ) = 64K (1 segment).  
Bit 12 is page invalid - paged out (PIC 11).  
Also holds Extended storage address - holds the two extra bits needed to go above 16 meg. (+ since CP itself does not use page tables it cannot go above 16 meg).

Virtual addr plus CR 1 take us into the segment table. There we pick up the page table address. We index into this using the page in the virtual address + find the real address of the page (or perhaps an invalid bit signaling that the page is out on DASD).

Prefixing - special translation for addresses between 0 and FFF (page 0). If, after translation, the frame number is 0, then the correct frame is substituted. CP uses the two highest frames in storage for RSAs.

CP can use real absolute page zero. Then, backwards translation occurs. If the frame specified is one of the high pages, then it get changed to zero!

SIO is slow. The processor stops until the channel has a path to the device. SIOF immediately results in a CC of 0.

CP runs disabled. Virtual Machines run enabled.

CP runs in supervisor state. Virtual Machines run in problem state.

CP runs DAT off. Virtual machines run DAT on.

Three units of work in CP:

1. Important CP (IOBLOK + TRQBLOK)
2. CP work that can wait (CPEXBLOK + CPEX)
3. users (+ the machines is enabled).

This is also the order of importance. The only reason to run a user is to be enabled so that an interrupt can come in and CP can get back to work.

Important work is IOBLOKs + TRQBLOKs. Each has a IRA (Int. Return Address) that is a LPSW to show how to get the user going again.

### Sysgen Preparation

Claudia

1. Find out exactly what devices exist - including all options (controller plugging, for example).
2. Special requirements like V=R?
3. Microcode assists? PMA, VMFA
4. Processor configuration: UP, AP, MP
5. Mode of operation: SPM?
6. What software will run under VM?

VS1 or VSE handshaking?

OS Compilers (must be ordered separately).

3705? VM only supports EP, no NCP or SNA.

OEM software - like GDDM (coded as saved systems in SNT)

7. Understand VM restrictions:

sys Planning Guide - Appendix D

Alt Path or Res/Rel?

Shared DASD between MVS & VM

8. Paging volumes adequate?

9. Will CMS have shared segments.

Plan DASD space allocations.

System definition files:

Directory (two copies) - lots of work.

DMKRIO

SYS - defines paging, nucleus, operator, etc.

SNT - ... saved systems & EP

BOX - logo

{ FCB  
UCS } Printer controls. Need only be  
{ DMLUCB } changed for special printers.  
{ DMKUCC }

PIA

PTB

Format disk space:

IPU FMT

### Service Programs

Service tape files:

1. ICKDSF - analyze + re-assign if necessary,
2. CP Format/Allocate Pgm (defines in 4K blocks)
3. DASD Dump Restore Pgm

CP Disk cylinder  $\emptyset$  must be specially formatted.

IPL DSF out of the reader. DSF is not used to format under VM. It is only used to check for bad spots.

Analyze - hardware problems

Init - volser

Inspect - surface check

MAPACT - only for FBA devices

REFORMAT - replaces volume label info.

<sup>vol 2</sup>  
INIT UNITADDRESS(cuu) <sup>150 / 160</sup> NOVERIFY VOLID(serial) <sup>→ don't check current label</sup> - <sup>new volser</sup> NEWVM  
optional → OWNERID(owner) VTOC(cyl, trk, extent) - (0, 1, 1)  
→ MIMIC(MINI(cyls)) DEVICETYPE(type)  
    memic 40 cyls,  
    not a whole volume

INSPECT UNITADDRESS(cuu) VERIFY(serial, owner) -  
NOCHECK ASSIGN TRACKS(cccc, hhhh) -  
PRESERVE

On S disk is IPLable version of DSF.

XEDIT INIT JOB A

GET IPL DSF S2

Format/Allocate must be run against any CP owned volume - Spool, paging, directory, sysres, etc.  
cyl 0, Trk 0, Rec 4 has bit map indicating use of device

01 Permanent Space (Non-CP use, Mdisks)

The rest of 150 must be allocated as perm.



- Ø4 - inactive Directory
- ØC - active directory
- Only the Directory can live at cyl Ø.
- Ø2 TDSK - temporary disk.
- Ø8 - page (preferred only)
- ØØ - Temporary - spooling + non-preferred  
paging
- 1Ø - Dump - 200 cyls.
- OVRD - override file. Permits the class of a command to be changed.

Put the nucleus at the tail end of the pack.  
(Only requires about 3 cyls).

### Dasd Dump Restore

Copy CP NUC from tape to disk.

DDR files from release to release are not always compatible. Therefore, copy DDR onto the tape ahead of the DDRed file.

The CP + CMS DDRs are not necessarily the same.

RESTORE NUCLEUS to get starter system.

### Second Level Sysgens

Planning Guide has the devices defined in the starter system.

Volser of second level disk must be the same as the CMS label of the minidisk.

EDVMIØ2 - userid  
TURKEY password  
CCCVME

VM screen - D PVM  
TSO - EC

### Directory

11-18-86

The directory defines the virtual machines. Claudia

Always located on sysres - a CP owned volume +  
pointed to from V014 label (Cyl Ø, Trk Ø, Rca 3).

The Direct file may be written out with DMKDIR  
using CMS DIRECT command, or standalone -  
IPL DIR. (SP PUN CONT, PUN Pgm, PUN Directory, IPC C).  
Use the CMS command.

The DMKDIR pgm only checks syntax, not  
minidisk overlaps.

\* in column 1 is comment line.

No continuation characters required.

Maint must have access to sysres to write out.

Spool + Page space can be allocated to a  
nolog userid + then Maint can link + do  
work if required.

The options Realtimer + ECMODE are necessary  
to IPL second level.

ACCT option - collect accounting info for machine.

## DMKSYS, DMKRIO, DMKSNT

DMKSYS - describes sysvas.

Describes all CP volumes, operator console, storage size, page + spool, monitor options, password suppression option, accounting parameters, system id (lower right of screen), + ...

SYSDOWN - must be 1st. Specifies up to 155 CP-owned volumes by VOLSER.

SYSDRES - 2nd macro.

volser, address (existence of alternate volume), devtype, starting loc of NUC in cyl - place at the end (3 cyls);  
All allocated out of perm.

SYSOR = real memory specification

sysmon - monitor spec:

designated userid to receive spool file class

Auto start?

Enable for classes (Perf, User, Dastap)

start + stop times

Record limit

Number of Buffers (3 or 4)

syslocs - must be the last macro.

Put paging on fast devices, spool on slow.  
Separate on volumes + channels. Isolate TDISK  
from anything important.

Consider multiple CMS system + Y disks.

HPO DMKSYS differences:

Directory has PMA option (to follow V=R option).

DMKRIO - HPO supports 32 channels 00-IF.

DMKSYS -

SYPAG replaces SYSORD. Must go after  
SYSOWN + before SYSRES.

Swapping must be specified first.

SYSCOR - allows up to 64 meg.

Can have multiple syspage macros - but will  
round-robin through the first set first!

DMKRIO - Describes I/O farm. VM will do a  
TIO to each device.

For 308x processors an IOCP is also required.

All macros must be defined in the proper sequence.

Macros:

Cluster - defines control units for remote 3270s.

Can be a 3705, 3274

Terminal - define 3278s + 3279s before  
3276s, 3277s + 3275s. This is required to  
use the new features of the new  
terminals.

Tony Chestang - friend of Phil, Mike, & Jill

RCTUNIT

RIOGEN macro must come last.

VMFASM is used to assemble these decks.

VMFASM filename ctrlfile (options

Use the same ctrlfile for RIO, SNT, & SYS.

Options -

DMKSNT -

To identify saved systems. Only has 3 macros. Defines what system & where on disk it is saved.

SAVESYS command provides name.

SNT must reside on a CP-owned volume, usually sysres.

— skip to next page —

## CMS Nucleus Generation

CMS will go on the 290 disk in lab. CMS lives on a minidisk.

CP is simpler than CMS.

CMS on minidisk has:

1. system files
2. CMS NUCs

Format 290 B

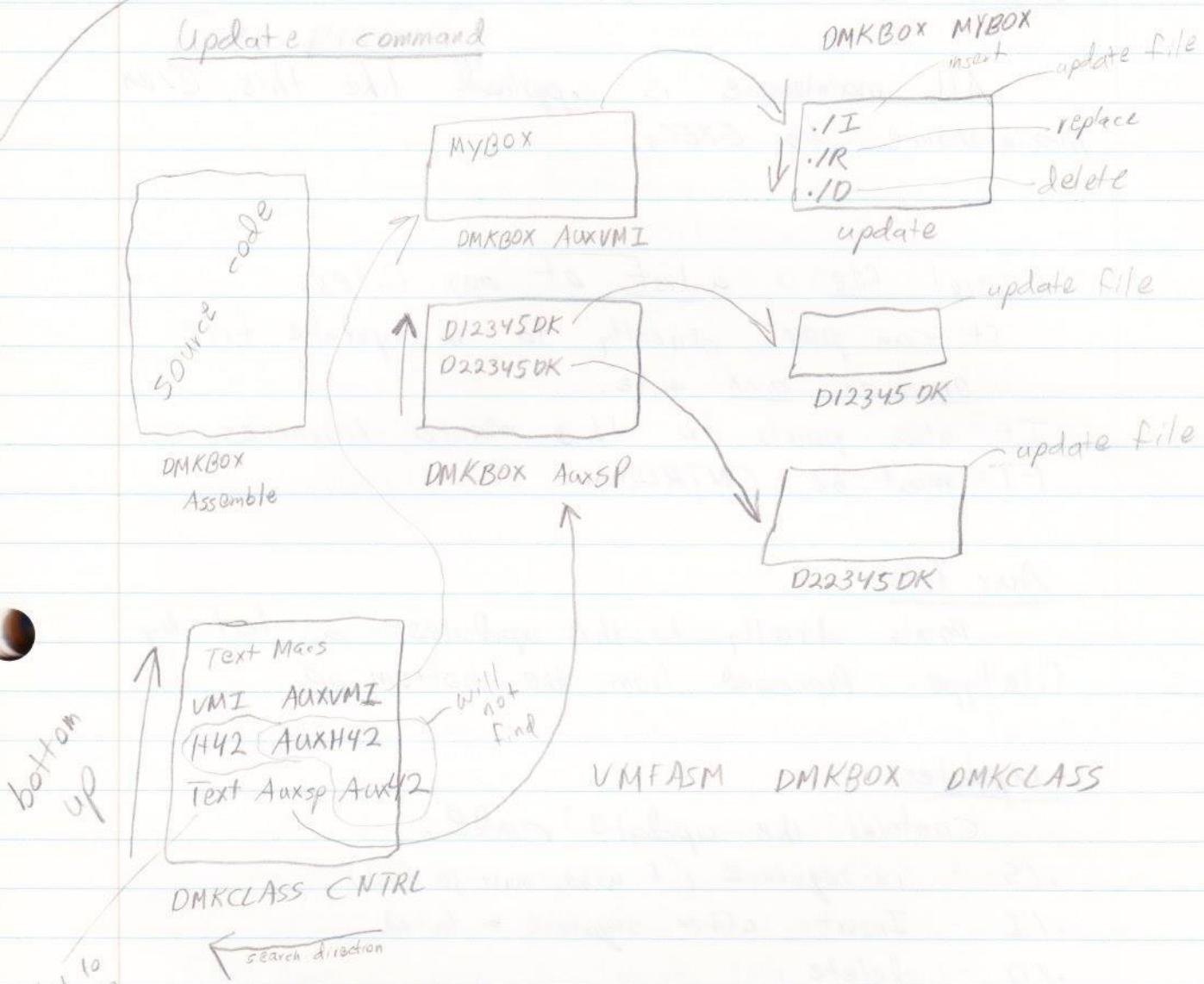
Format 290 B 1 (Recomp

First cyl is CMS file system & rest is for NUCs.

QV 290 will show 5 cyls; Q Disk shows 1 cyl!

CP Nucleus Generation

Update command



what to call the output

CNTRL file controls what updates go on the file. If an aux file is found, it is read. Each entry will point to a file of IEBUPDTE-like commands which will insert, replace, + delete code in the Assemble file. Each update listed in the Aux file is read. Aux file is read bottom to top.

Each line in the CNTRL file is read right to left. Once any aux file listed on a line is found, the other entries on that line are ignored. When the updates provided

by the auxfile are complete, the next higher line in the CNTRL file is examined.

All maintenance is applied like this, even maintenance to EXECs.

Control file is a list of aux files.

It can point directly to an update file or an aux file.

It also points to the macro libraries.  
FT must be CNTRL.

### Aux file

Points directly to the updates. - a list by filetype. Processed from the bottom up.

### Updates

Contains the update code.

- .IS - re-sequence (if used, must go first)
- .II - Insert after sequence # listed
- .ID - delete
- .IR - replace
- ./\* - comment

## CMS Nucleus Generation

[see previous]



11/20/85

Claudia

## Maintenance

PUT tape format:

File 1 - VMSE RV

2 - Memo

3 - List of pgm prod maintenance on tape

The next file will have an EXEC to load down the maintenance for each product

CP Aux

CP update

CP Aux for macros

Updates for macros

Source

CP Text

Do a VMFPC2 SCAN of tape file to see what's there.

All PUTs are cumulative from the base for a particular release (except for level set - which basically causes a new base tape to be released).

Order the PUT Bucket when ready to apply PUT. It may be wise to skip a PUT level occasionally.

To remove a PE comment out the Aux file entry & then re-run. (Check regen section of documentation).

L \* fixname to see which modules are effected by the bad fix. Copy over the IBM aux file & the comment the entry. Place the newly named aux file ahead of the IBM aux file in the search order.

## Local Mods

Copy the IBM aux file + give it a new name:  
AUXxxx

## CMS

### EREP

Classes C, E, + F are necessary to access the error cylinders. Only class F can clear the error cylinders.

### IPCS

Used to move dump spool files to disk for diagnostics.  
IPCS/E comes with SP4.

If you add maintenance to IPCS then you must use UMFASM. Otherwise, treat it like a CMS command:

Load pgm

GENMOD module

Copy module to CMS system disk

DUMPSCAN n ~ to look at dump.

IPCS needs the NUCMAP available to find entry points.

PROB - create/update<sup>a</sup> problem report manually in the absence of a dump.

APAR - prints a form describing a problem.

Amdahl Analyze is free to Amdahl customers.  
Its has better control block formatting features than IPCS.

## CMSBATCH

disk 195 is A disk. Some shops will create  
to CMSBATCHs - one w/ 30 cpls + one w/ 150 cpls.  
195 disk is erased between every job!

Control cards:

1JOB userid account ← no validity checking!

1SET time - seconds #lines #cards - limits

1\* - end of job

and any CP + CMS command.

CMSBATCH should be class G.

Clear out anybody else's LINKs prior to  
starting your job.

## RSCS

RSCS Direct - defines network.

Installation:

1. tape to disk 195
2. xedit RSCS DIRECT
3. " Profile
4. IPL. 195

RSCS Direct - do not put additional info in  
this file. Connections must exist.

RSCS is BC mode:

chan. 0 - Byte

1 - 5 - selectors

Can only have six channels.

RSCS maintenance files are kept on Maint's  
294, 394, + 194

Backup / Restore

DDR may be used to preserve the CP-owned areas on disk volumes. The DDR virtual machine must have MDISK statements for the entire pack. Remember to develop standardized restore procedures.

DDR may run under VM or standalone. Test the standalone before you need it. DDR exists in two versions, standalone + subsystem - thus the need to test.

For user areas, dump the whole volume via a full pack minidisk. To restore a minidisk, specify restore with only the cylinder addresses of the effected user. There are several products available to assist with VM backups.

Diagnose SVC 202

Opcode = X'83'

DMKHVC is a table of diagnose functions - a branch table. A user-written function would go into this table. 4 bytes:

X'83'	RX	RY	X'Function Code
1		2	3

0000 - 00FC Function codes are IBM pre-defined.

0100 - 01FC User definable function codes.

The ability to execute a given function depends on the command authorization of the virtual machine.

Use of diagnose is documented in the system programmers Guide.

## SVC 202

Used to call CMS modules. Can issue a CMS command from within a program.

## Programmable Operator Facility

Jim Turner

First fully implemented in SP3. The purpose is to compensate for the limited console support in VM. The heart of POF is a routing table which can intercept messages and let a specific program handle the function - which can be written in assembler or REXX. Messages may also be ignored - like logons and logoffs. IBM supplies several default options. The user can write pgms to, for example, att a tape drive to a user. This is a function which otherwise would be handled by an operator.

The messages may also be routed by POF pgms to another VM via RSCS. Thus, a remote site may be unattended. POF uses IUCV. The handling pgm must reside in a loadlib if in assembler. REXX will actually run faster than an assembler pgm.

## Monitor

System Macro in OMKSYS. Requires a gen to change.  
MONITOR command can over-ride defaults.  
To enable RESP requires a mod + re-gen. It comes disabled.

## Performance

Use INDicate and QUERY. The results of these commands vary with the authorization class of the issuing virtual machine.

SET QDROP OFF Q1 - forces user to stay in Q1.  
Use for VTAM virtual machine.

For CMS, never use a blksize greater than 12000. Otherwise, CMS splits it + retrieves two blocks.

With SP, keep page + spool close to the center of the volume. Put  $\frac{3}{4}$  on one side and  $\frac{1}{4}$  on the other. Then, when half full, it is balanced. HPO is smarter and this is not necessary.

Get CMS wac off of sysres. Put 19D + 19E on separate channels.

In an application development environment, separate the 191 disks of group members since they may tend to all be busy at the same time.

## Canned Procedures

Many EXECs exist to aid in the installation of VM. User must access the right disks at the right places & then an EXEC will format, download, update, assemble, ask a few questions, punch the deck, & IPL.

# MDF

12-17-87  
Gloria

4K - storage can be dynamically shifted.  
Also, we can have dynamic partition & join.  
↳ requires parmlib change to prevent long-term page fixing. May be useful for raiding a test domain for its storage during times of need in the production domain.

## CPU Allocation

Domain	# CP * Target %	= % allocated
1		
2		
3		
4		
sum		< 100 * #physical CPUs

Make a chart to schedule multiple proc / mult domain

	CPU 1	CPU 2
Dom 1		
2		
3		
4		
Total	100%	100%

← this must be 100% (or less) for each CPU.

The scheduling parameter on SS is global & permanent (written to 1SYS).



Have all pages face the same way.  
Get rid of Async msg on foil examples.

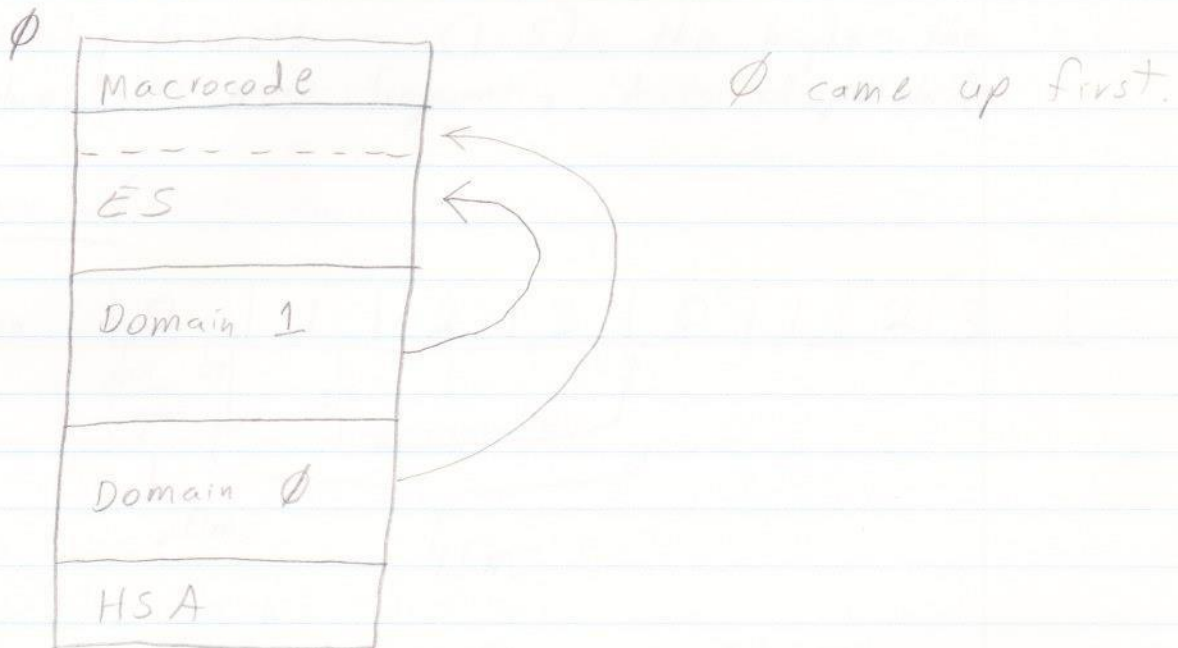
Sheet of equations.

Format of MRSD data - DSECT

Steve guarantees that a 5890 console  
can run on a block chan.

A one LP system without CPU affinity will still come down in event of a MCH check.

### Storage



There is no  
AMD IOCP 6.0C  
This must be  
EREK.

AMDIOP 6.0C now supports  
classic 580s.

5890 standalone now supports 3480s.

A user determined CPID (on CX frame) has  
replaced the old LP#

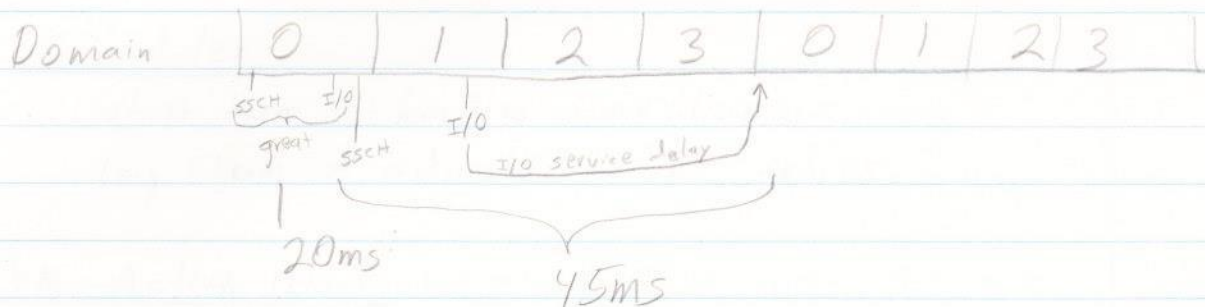
## Performance

CPU overhead has 2 components:

1. Domain switching.
2. Buffer pollution.

Scheduling Parameter - (1-5) - the higher the value, the more frequently it is dispatched.

### Low SP:



High SP will reduce I/O service delay:

0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 |

Domain comes active more frequently + picks up these queued I/O rpts.

ROT → use higher SP for interactive domains.

Due to I/O service delay, MDF may nullify the value of EDAS or cache.

Shared DASD may become dominated by a non-MDF system. Also, reserve/release becomes elongated.

A logically partitioned system has a 500ms time slice (this is the same as for a non-MDF Amdahl machine).

### Scheduler

short term - handles timeslices.

long term - adjusts target values.

### VM Active Wait

PSW      TIE kkkk CMWP  
          011 3 1000  
          ↑ bit 5

If one domain is a CPU hog (batch...), protect the interactive domain by giving it a minimum. Else, time will migrate to the hog.

ROT: for beginners -

SP = 3

target = max

Change gradually.

## RMF Measurement

clocks:

TOD - wall clock

Clock Comparator - causes Int when  $<$  TOD

CPU Timer - decremented only when CPU is in operating or load state.

Interval Timer (370 mode) - used by TPF (ACP)

When a domain is not dispatched, the effect is the same as if the stop button had been pressed.

TMB219 - Customer version of explaining RMF wait reports.

Wait time is accurate. The sums of waits for all domains is the complex wait. What we don't know is how to attribute the busy among domains. This works even with floating percentages.  $100 - \text{complex wait} = \text{complex busy}$ .

With fixed % we can get domain busy:

Fixed Domain % - % domain wait = % domain busy of complex

1. How busy is entire 580?

$$580 \text{ wait \%} = \frac{\text{sum of measured waits}}{\text{interval} * \# \text{ physical processors}} * 100$$

2. How Busy is each domain?

DORT = Domain Residency Time - time domain is dispatched on 580.

$$\text{DORT} = \text{interval} * \% \text{ target (for static \%)}$$

$$\text{Domain wait \%} = \frac{\text{DORT}}{\text{interval} * \# \text{ phy proc}}$$

3. What portion of the 580 is each domain using.

Total 580 Busy time =

$$(\text{interval} * \# \text{ phy proc}) - \text{sum of waits}$$

## MDFWatch

Shows how busy the domain is relative to the complex.

MRSD can show busy relative to the domain - but not realtime. This must come in through an exit.

# VM HPO CP Internals

2-29-88

Hellis Johnson

EDHPIxx - userid

SP PRT RSCS

Tag Dev pit complex2

EPSPT1

} to print

191 - A disk

194 - CP text

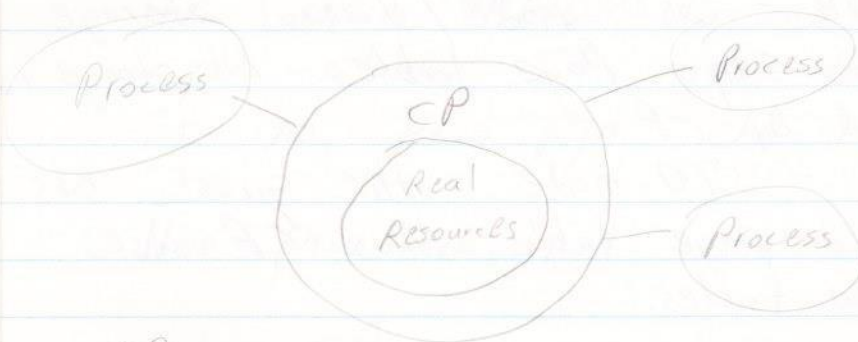
19D - Help

294 - CP updates

321 - class

322 - class

394 - CP Source



CP

Manages Real Resources (Mem, I/O, CPU)

Schedules

Provides services (with a degree of device independence) DIAG

CP simulates a console:

Begin → start button

CP Read — stop

Ent — Int

Ind — Meter

sys clear — reset

IPL — Load

Adstop — address

Trace → single step



CP provides 16 meg. of storage. A V=R user can get more.

User storage is furnished via segment & page tables of 370 architecture. Guest virtual storage is supported via shadow tables.

Seg & Page tables are in high tree. DAT utilizes these tables, they are a part of the hardware architecture.

Swap Tables are not part of the arch; these point to pages in aux storage.

### Virtual Machine

BC mode - 360 mode / Virtual Storage  
Uses Seg & Page tables. All storage is managed by CP

EC mode - 370 mode: the guest has seg & page table. Thus CP uses shadow tables.

BC mode addressing:

VMBLOCK  $\rightarrow$  seg table  $\rightarrow$  Page table  $\rightarrow$  Swap table  
 $\rightarrow$  CR1 (STOR)

EC mode addressing: shadow tables are required. CP can handle up to 12 shadow tables at one time.

Substantially more overhead to run an EC virt machine.

Devices:

RDR PRT, PUN - totally simulated

DASD - can be real; are shared

MDISK statement

other - must be attached as real; tape, for example cannot be shared. If a real prt is ATTACHED, it cannot be shared.

Instruction Processing - CP intercepts instructions which effect the whole hardware. Virtual machines run in Prob. state! The machines believe they are in Supervisor state  $\rightarrow$  SIO, CPSW, etc.

The real PSW is prob. state, causing pgm chk to trans control to CP which then checks VMBLOCK to see if in virtual super state. If yes, then simulated.

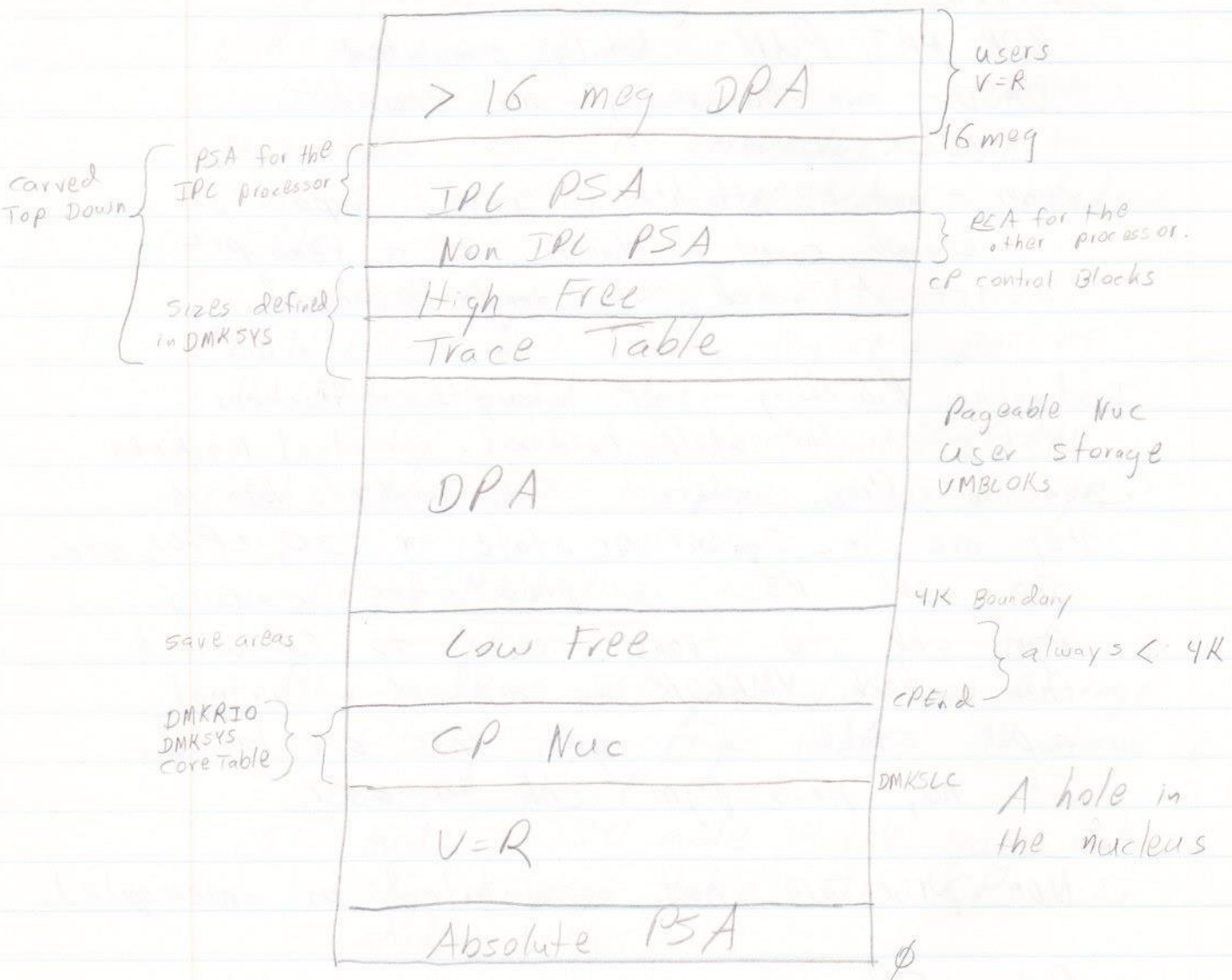
If no, pass pgm chk to user.

Non-priv are not simulated or intercepted.

Only CP can run in supervisor state.

Timeslice is protected; user can never really disable for ext interrupts.

# Storage



If high free fills, CP will take a page from DPA (page Extends).

If the >16 meg area is used by V=R machine, then all storage between V=R area + upper area appears as offline to the guest.

Insure core table is just the right size.

CP names:

DMK<sub>xxx</sub>yy

→ use of module

RIO - yy = CH, CU, or DV

SY5

yy = entry point

SCH

DSP

VAT - builds shadow tables

SCN - scan; generalized search

yy = RU or VU

Regs:

0 - length passing (Free/Fret, QCN, etc.)

1 - Free area adr; buffer pointer

2 - Cntl parameter; PTR, QCN

I/O { 6 - real/virt, chan pointer

7 - " " CU

8 - " " Dev

10 - stackable entity - adds to 9

11 - VMBlok of cur user

12 - first base reg

13 - save area pointer (18 fullwords, Pointer is to last 16)

14 - return link pointer

15 - Goto pointer

CP SVCs: used by CP, not users.

0 - abend

4 - reserved (abend)

8 - link request } builds savearea & branches

C - return }

10 - release savearea

14 - obtain savearea

18 - AP/MP; SIGP

## CP Linkage:

BALR - uses BALRSAVE in PSA (only one)

GOTO<sub>macro</sub> go to dispatcher

implied exit; uses svc 8

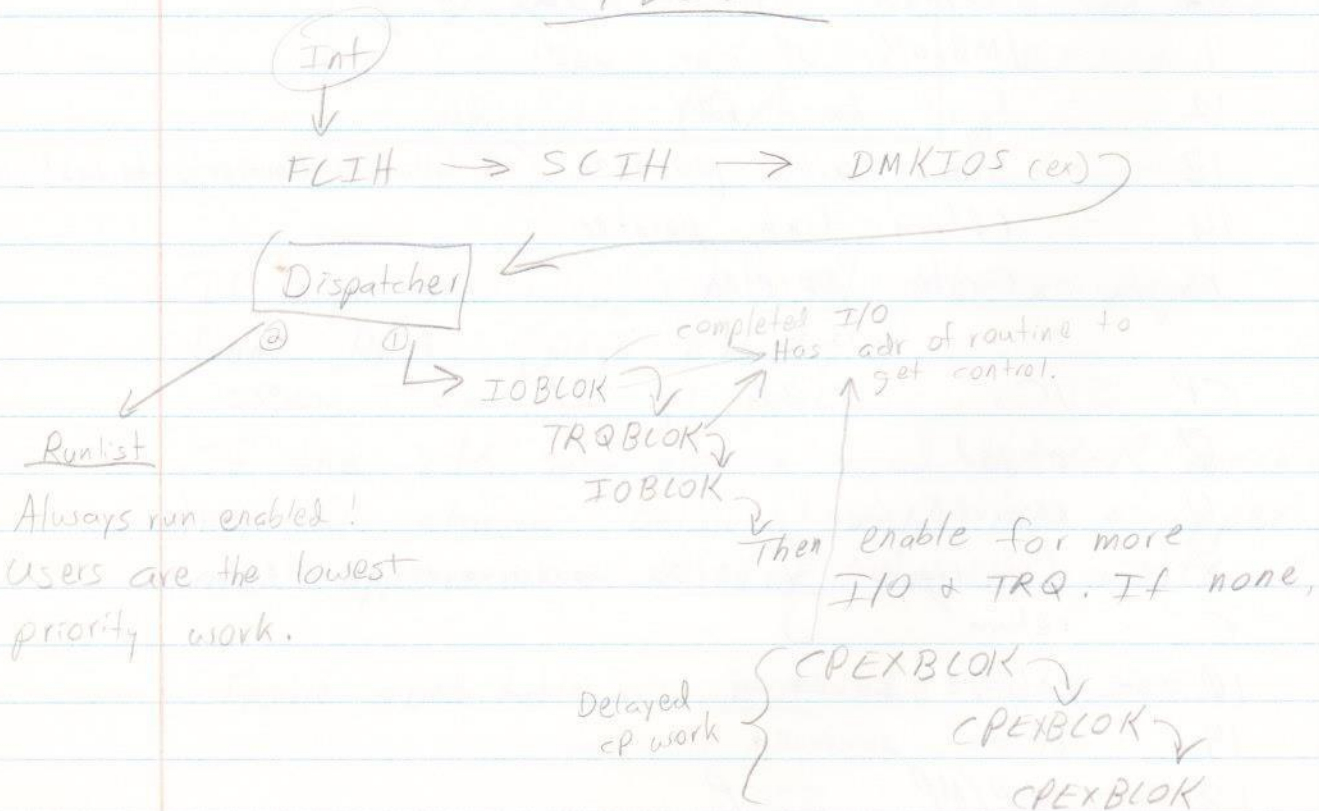
SVC Linkage - DMKSVC gets savearea

The CALL macro decides whether to BALR or SVC. DMKFree uses Freesave as a private holding area. Even routines BALRed to, can get free. CALL macro has a list of routines which can accept BALR. Most transfers are via SVC.

BALR modules must be CP resident.

All interrupts pass control to CP code.

## FLOW



## Architecture

EC PSW - status

next inst. adr. (IS)

Bit.

status

15

All users are in prob. state

12

EC/BC - CMS runs BC mode in virtual PSW.

SETECMODE ON does not force EC mode unless guest loads an EC PSW. This allows an ECBLOCK to be created. In CMS, permits use of chans over 5.

5

Trans/Non Trans - Trans means IS is a virtual adr. - Used by all users.

Non-trans - CP; Real adr

6+7

EN/DIS - I/O, Ext en for virt. machines; Disa for CP

Key - 4 bits

Pgm Masks - Set in real PSW by user.

20 - Fixed point overflow

21 - Decimal

22 - Exponent underflow

23 - significant

070D

- typical for user

0002

- Disa wait

BC PSW - bit 12 =  $\emptyset$

No DAT bit

Channel masks are in PSW

### Interrupts

Allows CPU to change its state via PSW swap.

Type code

Six kinds:

Old PSW

Ext -

X'18'

I/O -

X'38'

Mch -

X'30'

Pgm - cannot be masked

X'28'

Svc - " " "

X'20'

Restart -

X'8'

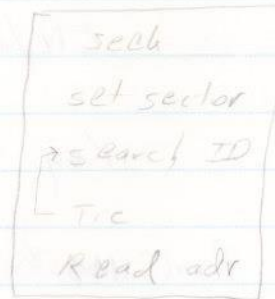
### 370 I/O

Build Chan Pgm  $\longrightarrow$

CAW points to CCWs

Put dev adr in Reg

SIO of Reg of Dev



check keys with CKD

LA R2, CCWs

ST R2, CAW

L R1 $\emptyset$ , Devaddr

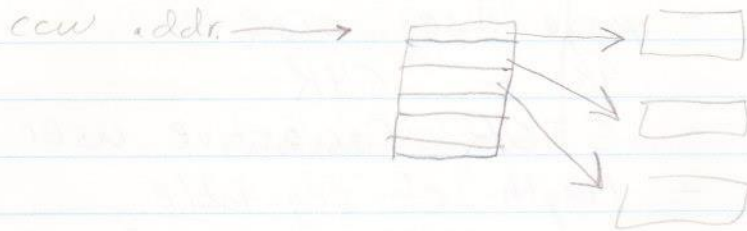
SIO  $\emptyset$ (R1 $\emptyset$ )

### I/O Int.

CSW points to CCW after failing CCW

IDAW -

ccw addr area  $\rightarrow$  table of addr  $\rightarrow$  stor.



HPO uses SIOF.

Timing - 4 clocks:

TOD - beginning May 11, 1971

Doubleword

Bit 31 increments gpx 1 sec.

Bit 51 is microsecond.

Runs always

Clock Comparator -

Doubleword

When = to TOD, ext int. 1004

CPU Timer

Doubleword

Decrements

Causes Ext Int 1005 when negative

Does not run in stopped state

Used when setting timeslice

Interval

Full word

Located in Mainstore

Old 380 Timer!

Bit 23 decrement 300/sec

Decrements

Causes Ext Int 0080 when negative

Does not run in stopped state

Used to set Dispatcher Timeslice



DAT - 5 parts:

PSW - EC mode bit (12) on  
- Translation bit (5) on

CR0 - page + seg size  
4K + 64K

CR1 - STOR for active user (VMSEG)  
- length of seg table

Segment Table - 4 byte entry for each  
64 K.

Page Table - available  
- real page frame address

CR0 - bits 8-12 = 10000

to indicate 4K pages + 64K segments

CR1 - bits 8-31: seg table origin.

bits 0-8: add 1 + multiply by 64

size of table      document this much storage

0 = 64 bytes = 1 MB

3 = 256 = 4 MB

F = 1024 = 16 MB

Seg table Entry - 4 bytes

when calculating page table addr, consider  
bits 29-31 to be 0.

Page Table Entry - 2 bytes

maps 4K. Hold Real Frame Address.

Seg # is first 2 nibbles

3rd digit is page

01 = 1  
10 = 2  
11 = 3

Real Frame Addr | 1 | 0 0 |

12  
page invalid  
Extended Storage

virt addr = 555764

go to '55' into segment table + find  
F0764762 → PT

→ Get 5<sup>th</sup> entry in PT = 444 page number  
Displacement into page is 764, therefore,  
real addr = 444764

If Extended Storage bits are on, move  
to front of address. If 11<sub>2</sub>, then,  
real addr = 3444764

Highest addr is 3FFFFFF = 64 meg.  
This is the 64 meg addressing limit.

If AP or MP, prefix Reg specifies page  
which holds PSA. Format = xx FFF xxx page #  
After real address is obtained, it is compared  
to prefix reg. If equal, real adr is changed  
to 000. If real adr page = 0, then change  
it to prefix reg page #.

# Control Block Overview

## PSA

- Anchors everywhere
- counters
- save areas (BAKRSave, FreeSave)
- CP info
- CP status (AP, MP, UP gen + running)

### → Coretable

- 1 entry per page frame.  
who owns, etc.

### → PSA Extension (1 per processor)

Dispatch List → IOBLOCKS, TRQBLOCKS, CPENBLOCKS

Run List → VMBLOCKS (started to run + want to continue running)

True Run List → Immediately ready to run. NOT waiting on I/O, Page...

### → Real I/O Blocks

Built from DMKRIO; CH, CU, DEV  
↳ Device Dependent Blocks

### → VM Bloks - chain of all user VMBLOCKS

- ↳ Segment, Page, + Swap tables
- ↳ Virt Arch Blocks
- ↳ Virt I/O Blocks
- ↳ Virt Spooling Blocks

### → Real Spooling Block

- PSA - hold pool of useful constants  $\rightarrow \emptyset, 1, -1$
- trace table pointer (since both CPUs write to same table, the absolute PSA keeps this).
  - Save areas
  - Pointers to resident routines (for many BALRs)

PXA - came with 3.4  
 1 / CPU  
 contains Dispatcher Queues.  
 Subpool info -  
 CP storage pools of various sizes

### Real Machine

portable

SCBLOK - Swap C.B. Built when users are logically swapped. Anchors pages that will be swapped when user is physically out

SSBLOK - Swap Set Blok - where the swap set went to aux.

RCHBLOK }  
 RCWBLOK } I/O  
 RDEVBLOK }

### DASD

DASDALLOC - cyl  $\emptyset$ , rec  $\emptyset$ , hd  $\emptyset$

Reads allocation map of volume

AREABLOK - defines use of each rec out on DASD. 3350 - 120 pages/cyl  
 3380 - 150 pages/cyl

Is page in use...

sysplist - Syspage macro info. specifies hierarchy: swapset + preferred paging

SFBLOK - Spool File Block

4.2 + below - high free (9900 limit)

5 - separate address space.

Real I/O Blocks

Built contiguously in memory.

PSA - ARIOCT - chan index table

32 halfwords containing index into chan table.

FFFF = undefined

ARIOCH - start of chan table.

Add index to this to get entry of particular chan.

RCHBLOK -

RCHCUTBL - cu table.

32 entries / # CUs on chan.

Get index from ARIOCU in PSA.

RCUBLOK

RCUDVTBL - Dev table

can specify two chans (RCUCHA, +B) - backward pointers.

RDEVBLOK -

can specify two control units (RDEVCUA + B), backward pointers.

IOBLOK - chained off RDEVBLOK (@ RDEVFIQB. They may be chained off the highest busy component.

Locate I/O Block:

ARIOCT - chan index table ↑

ARIOCH - start of RCHBLOKS

CU - RCUBlokS

DV - RDEV

Chan index table - halfwords index

$ARIOCH + \text{index} = RCHBLOK$

CU index table - 32 halfwords index

$ARIOCU + \text{index} = RCUBLOK$

Dev index table: 16 halfwords index

index compressed

$ARIODV + (\text{index} * \text{compression}) = RDEVBLK$

Compression

4 bit shift

Device Dependent

UR - RSPCTL  
RSPXBLOK

Term - CONTASK

Transmission CU<sub>s</sub> - NICBLK (network interface - defines 3705s)

BSCBLOK - bisync devices

VM BLOK - one per logged on user.

- alloc in DPA on 128 byte boundary
- chained in cyclic chain

(QN starts with your VM BLOK)

- Virt Machine Info

- runnable
- anchors
- virt regs
- machine stats
- virt PSW

Virtual PSW - save old PSW into virtual page 0 + interrupt code. New PSW is placed in VMPSW in VM BLOK to prep for next dispatch of user.

## Storage

Segtable

Pageable

SWPTABLE - paging

SCBLOK } swapping

SSBLOK }

## Arch

MicBLOK - what microcode enabled for user

ECBLOK -

STOBLOK - multiple STORs indicator

XINTBLOK - ext int block for saving ext int for user; will be presented when user enables.

## Virt I/O

Same arrangement as real, but anchored in VMBLOCK.

## Scheduling + Dispatching

1. IOBLOCKs + TRQBLOKs, intermixed
2. CPEXBLOK
3. VMBLOKs - backward/forward chained.

Run List

True Run List

## IOBLOCK

addr

forward/backward pointer

CSW (to give to user)

IOBIRA - routine to get control.

Built at I/O request time + posted to a dev. At I/O interrupt, CP places it on a dispatcher queue.

TRQBLOK - bit 0 is 1 to distinguish

it from an IOBLOCK. This is why May 11, 1971 is the oldest date!

- CP + user TRQs are on same queue.
- Address of routine to get control - TRQBIRA
- Pointer to VMBLOK of requestor

CPEXBLOK - delayed work

- contains Regs

- forward/backward pointer

- address to get control

- can be chained from IOBLOCK. When IOBLOCK gets dispatched, it adds this to the queue.



# Real Storage Management

3-1-88

Jim Turner

STOR must be on 1 meg boundary. Therefore if DET STOR is 512K, VMSEG will still be 1 meg, but VMSTOR will describe exactly how much can be used.

Page - 4096 bytes of data

Frame - in mem (last 12 bits are 0)

Slot - on aux storage

True Working Set - pages ref while in queue. These will be logically swapped out.

Trim Pages - all pages not referenced while in.

Swap Set - group of ref. page - True Working Set subdivision. Trim pages will be added to achieve the Excess over multiple SS size are placed on trim.

Logical Swap - still in store, on swap g.

Phy Swap - swap sets moved out to aux.

Swap Area - SW on syspag macro. System will swap even if no SW defined.

Interactive Buffer - Portion of DPA reserved for Q1 users. This is really an algorithm that effects eligible list candidacy. This # pages is subtracted from available pages when CP decides whether to run a non-Q1 user - more likely to go to Elist. IB is set too high by default - 57%  
SET SRM IBUF to 10%.

Q2, Q3, + slow Q1 users (SPT SAM  
ZBUFF BUFTIME xx) appear to be entering a  
smaller DPA than busy Q1 users.

DPA is what's left after NUC, V=R, + low free  
are defined bottom up, an PSAs, High Free,  
Prime Free, + Trace are define 16meg line downward.  
In V=R paging is controlled by the guest. PMA + VMPA  
permit MSB/SP to use a V=R area above 16 meg.  
Also, CP does not have to do CCW translation of  
the CCW addresses (V=R page zero addresses  
must still be translated).

Size of CP nuc may be adjusted with  
some limits. The Nuc resident coretable is  
rather large and cannot be moved out. Nuc  
ends at DMKCPEND. Above this, to the next  
page boundary, is low free. In SP,  
this is an SVC save area. In HPO it is  
treated as high free. SP benefits by having  
non-paging low free. Some sites artificially  
increase DMKCPEND to push it above the page  
boundary.

CP runs DAT off + therefore cannot use  
the 2 bits in the page table for addressing  
above 16 meg. Only the DAT on users can use  
this area.

High Free is defined by SVSCOR macro,  
FREE operand. Default is 3 page frames for  
1st 255, 1 page for each 64K to 16 meg, and  
1 page for each 256K above 16 meg.

Excludes  
V=R area

If AP/MP, then add 25%. The default is usually too low. Try 2 page frames per logged on user. If high free is short, it will EXTEND into DPA. This is bad. Tune for virtually no extends.

Acct records live in high free until enough accumulate (80 bytes each) to cut a spool file.

In HPO, prime free storage is used for IOBLOCKS, TRQ, SVC save areas. Default is 10%; should be 2%. This holds CNTL blks that are performance oriented. Check "Prime Misses" + reduce until some appear, then raise some.

Trace is usually too high! 4K is minimum + may be fine. Only SMART cares. Excess here goes to DPA. Trace cannot be change on SP w/o a mod.

UMBLOCKS live in High Free in SP, and are locked in DPA in HPO. These are placed on a 16 byte boundary.

### Storage Protection

The 2K key instructions (SSK, ISK, + RRB) can be used on post-3081 machines if CR0 bit 7 is set. CP checks this before dispatching. If not set, then pgm chk.

DMKSTA checks at initialization to see if the processor is 2K or 4K keys. This also checks for VMA ucode support.

PSA + 3CC points to CORETABLE (ACORETBL).  
16 byte entry per page frame - two formats:  
if on a list (free, flush, swap) the entry  
will have forward/backward pointers on the list.  
Else, adr of VMBLOCK (owner) + lock count  
(is CP using, is I/O using, then page is fixed).  
Coretable describes all storage. CP, V-R,  
etc., is permanently non-payable. Lock count  
will always be  $> \emptyset$ .

Core table is in DMKSYS.

In DMKSYS,

1. DMKSYSRV - real store size not including PMA.
2. DMKSYSRM - actual real storage size.
3. DMKSYSRC - real store size including PMA.

1 + 3 are set at SYSCOR Macro. This determines  
size of coretable

Entry: 5 types.

1. Freelist: CP wants users above the line.

- Below 16 meg - DMKPTRF1

- Above 16 meg - 2

The user page  $\emptyset$  is kept below 16 meg line  
because CP references this for virtual I/O.

I/O pages must also be below.

Freelist pages are immediately available.

CP maintains some pages on the freelist.

It will anticipate page outs to keep  
this up.

2. Flushlist - DMKPTRUI

This is a list of trim pages.  
Marked as invalid + maybe reclaimed.

3. Swaplist - logically swapped out pages.  
These have been referenced.

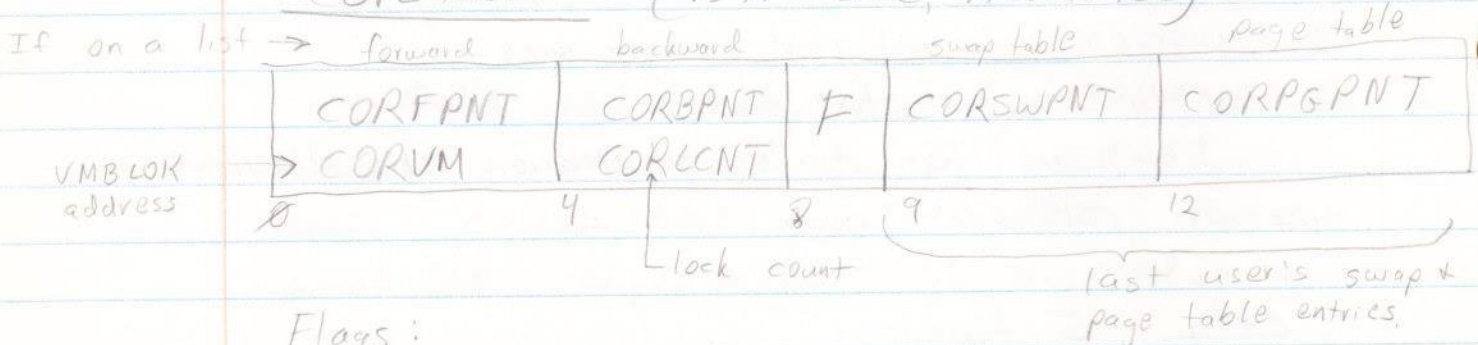
One swaplist/user. Chained  
from VMBLOCK → SCBLOCK → page.

4. Not on a list -

User page -

System -

Core table (PSA + 3CC, ACORETBL)



Flags:

- 80 - locked for I/O
- 40 - " " console
- 20 - flushlist
- 10 - free list
- 08 - shared
- 04 - swaplist (HPO only)
- 02 - CP
- 01 - Disabled or not available

If in transit, byte 8 holds CCW opcode.  
(06 read, 05 write).

Position in coretable determines real address of page frame. They are sequential.

$(RPN * 10_{16}) + ACORETBL = \text{coretable entry}$   
CP deals with Real Page Number + uses this to find the right entry.

Swap list + page tables are anchored in the user's VMBLOCK. CORETABLE points to the exact entry in each. Swaptable entry holds disk location of page.

Swaptable has 16 entries:

Flag	V	key	key	DASD adr
------	---	-----	-----	----------

↳ virtual page number

If coretable entry describes high Free, the first word in EDCCIC "FREE."

Or → "OL" if offline

or → "FREA" if cache aligned DPA (VMBLOCKS)

One 16 byte entry for each 4K. This can take a lot of space.

Lists: DMKPTR points to 1st entries.

Flushlist - trimset goes here (not referenced during last run. May be reclaimed, or used to replenish the freelist).

DMKPTRS - HPO - only CP module that runs DAT on to move pages above 16 meg.

Freelist - when user logs off, or IPL clear, then pages go immediately here (also DIAG 14). Else, pages stop @ flushlist. These pages get here voluntarily.  
Involuntary - from flushlist, swaptlist, or stolen.

Swaptlist - working set. 1/user. These are now logically swapped out. This is purgatory.

### Q-Drop

At end of timeslice the user loses working set. Runs through pagetable (to limit of VMSTOR) + invalidates (bit 12) every entry. CP picks up the real frame number + finds the coretable entry for each entry + places each page on a list. They go to either working set or trim set + then subsequently are placed on the flushlist or swaptlist.

### VMBLOCK

↳ SCBLOCK ↔ SCBLOCK ↔ SCBLOCK ↔ ...

↳ Trim set  
↳ Working set

SCBLOCKs have counts of pages in each set.

The SET Min Working Set Size commands will be applied after winnowing. Pages are borrowed from Trim to working to fill out. This size should be a multiple of the swap set size! Because swap set size adjustment comes next & the non-multiple set would be paged, not swapped - undoing the intent.

If the user is SMALLV then it gets to keep all pages in working set. SMALLV is user who used  $< 1/16$  Quantum. Quantum is timeslice allowed. This is indirectly based on the proc. speed. 5890 type is about 5ms/Quantum. User can't hit many pages in  $1/16$  Quantum.

Logical Swap - add user SCBLOK to chain of SCBLOK. DMKPTSIF + IB anchor interactive + non-interactive queues.

Page Flush - Trim set pages are moved to Flush list. Interactive user pages go to bottom; non-inter go to top. Pages are taken from the top to fill the freelist. DMKPTRUM points to top.

Freelist - replenished by DMKSELECT in HPO. Its goal is to keep the freelist full. Page + Swap outs occur only to replenish the freelist. Freelist size = 1 page per user in runlist + 2.



CP prepages 2 swap sets (default) prior to dispatching an interactive user. Be sure not to exhaust freelist by pre-paging too many swap sets. SET SRW Minnumss + prepage are closely related.

### Swap Out/Page Out

Garbage Collection - time dependent. Since CP doesn't Q-drop CP pages in DPA would never get to free list. Garbage collection scans  $\frac{1}{8}$  of DPA beginning around every 4 sec looking for unchg, unret pages. Entire DPA is scanned after 8 cycles. Uses RRBE to reset ret bit. If already off, the page is placed on freelist.

Average life of garbage page is 48 secs.

HPO 5 Garbage collector is enhanced. Rel. 5 does less at Q-Drop. It will catch more user pages.

---

DMKPTR F1 = 1st CTE on < 16mb Freelist

2 = " " " " > " " " "

FN = count of # pages on < 16mb

N2 = " " " " " > " " "

DMKPTSIF → Interactive logical swap queue anchor

BF → Non " " " " " "

CP tries to save I/O. It puts swap sets out before draining flushlist. Pages of users logically swapped for 100 seconds are moved out with one SIO & lots of pages are freed up. This beats 1 sio/page from flush list. Making the time shorter will increase swapping. This can be desirable.

Once old non-interactive users are gone, then drain flushlist. A vFault (virt fault) will look for a page above 16 meg. Else, page fault is for CP use & should be below 16 meg.

Then, swap out non-interactive users from logical swap users.

Then swap out interactive.

phase 1 Finally, coretable scan! Not locked, not CP, not on a queue, not changed, then invalidate & steal.

If a page below is needed & a page is available above 16, then simply move a page above & free the frame below.

DMKPGT - find the best spot on disk.

PAG - schedule page write.

coretable phase 2 - quiesce other proc & take everything including reserved pages.

## Replenishment Summary

1. Disposable Page Collector (Garbage Collection)
2. Old non-inter working sets
3. Flush list
4. Swap queues
5. Core Table scan

DMKPTRFN - # pages avail (Freelist)  $< 16m$

N2 -  $> 16$

SS - # page steals systemwide

FF - # pages moved from flushlist

FØ - # times freelist empty

PR - # pages reclaimed

CT - # calls to DMKPTR

SW - # pages being swapped  $< 16meg$

S2 - " " " "  $<$

ES - " " " " stolen from in-g users

NP - " " avail in DPA

Some of the counters are reset at Monitor interval.

## Free Storage Management

This is management of high Free - DMKFRE + FRT.  
FRE still serves as the entry point for FRT. (DMKFRET)

Types of request:

Normal -

Big - Frels list

Little - from subpool (of appropriate size). <sup>separate for each CPU.</sup>

Prime - set size blocks for performance oriented control blocks. Separate for each CPU.

Cache Aligned - build VMBLOCKS at logon.

DMKFRE - contains code to allocate storage + manage it.

DMKFRT - logic to deallocate.

PXA - subtable  $\rightarrow$  subpool table chain anchor.  
2, 4, 6 doubleword chains. These are dedicated to a CPU. Alloc is fast.

DMKFRE holds FRELS list. Holds table of large free areas.

FRENUM - 8 bytes past FRELS; a count of # blocks on FRELS.

PXA - anchor for prime storage blocks.

DMKFRE - anchors cache aligned blocks which are available.

DMKFREE - entry point:

LA R0, IOBSIZE      # of doublewords  
CALL DMKFREE

+ L R15, AFREE

+ BALR R14, R15

On return, R1 has adr of storage.

Requested size will be rounded up to match a subpool. Area returned is not cleared. Storage must be explicitly release when finished.

Small subpool request: (less than 1024 doublewords, 8K!) (SP- 33 doublewords) - Round up request, check subpool pointer in PXA. If subpool avail, move pointer to R1.

If none available, go to FREES list. Pointer is  $\emptyset$  if none available.

92 Subpool sizes! (SP has 11).

2 to 128 doublewords in increments of 2 bytes.

160 to 1024 " " " " 32 "

The table has two entries for each pool size; one non-DPA followed by one DPA, etc.

DPA + non-DPA subpools are chained separately. Also, each CPU has its own set. Non-DPA subpool blocks are timestamped.

Once/hour subpool merge tries to combined blocks to place on FREES. This also happens whenever a user logs off. If a pool block is contiguous with FREES then combine.

Subpool Entry -

Doubleword in length.

Address of Free Queue Element

# of elements on chain

FQE adr	# FQES	Mon. Seg #
4		6

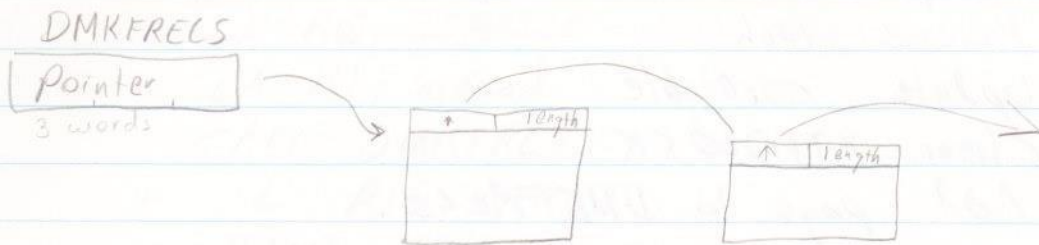
Large Block Request - DMKFRECS -  
entries are in ascending address sequence.

1st word = pointer to next

2nd - unused

3rd word - # of blocks

Each block has a length header. in 2nd word.  
First word is pointer to next block.



FRECS list is run, looking for first  
or last equal request.

If request is for greater than 28 doublewords,  
look for last equal entry. If no fit, take last  
larger size fit.

If no storage on FRECS, then EXTEND. Call  
DMKPTR to get DPA block + search again.

During Extend processing, no useful work is being  
done. Extend occurs one page at a time +  
if the request is for more than 4K,  
pages will be picked up until enough  
contiguous storage happens to appear! It  
could collect (page out) 60 page in order to  
deliver 5K for CP use.

## Extend Processing

Bump Extend Count.

Set XTNDLOCK (X'349' in absolute PSA) so dispatcher won't start anything new.

Decrement # avail pages

Get lock (DMKLOCK)

Call DMKPTR for a page frame - if a page out is required this is a performance killer.

Release lock

Update corotable

Clear XTNDLOCK

Add page to DMKFREES.

And think what happens if 4K is not enough.

## Prime

DMKFREEP - entry point. Prime is a permanently alloc pool of 16 Doublewords. 10% Default is too high. Change to 2% of high free.

If no prime is avail, use normal free storage subpools. No great loss.

Each proc. has own prime. Once/sec the pools are balanced between processors.

Cache Aligned Pool Request - aligned for HSB access. Six pool sizes but only 96 DW is being used (for VMBLOCKS). DMKFREC is entry point. Pointer is DMKFREA. If none avail, another page is taken from DPA + carved up. Like a permanent extend.

With 4.2, storage is released via DMKFRET, but entry is via FRE.

```
LA R0, JOBSIZE  
LA R1, areaadr  
CALL DMKFRET  
+ L R15, AFRET  
+ BCLR R14, R15
```

Subpool blocks are pushed on top of regular subpool chain. Prime go on top of PRIMHDR anchor.

If a whole block belongs to DPA, an unextend occurs.

Unextend occurs unless the routines are DMKQCO or DMKQCN which are known abusers. These guys loop on acquire + free. DMKFRET does not unextend, it holds block assuming that another request will be along shortly.



Storage Pool Merge - Once/hr or at every logoff, DMKFRET checks every subpool to see if it is adjacent to a FRELS area. If so they are merged.

Non-DPA blocks are returned to free storage if at least 60 seconds old, (for small blocks), or 15 secs for large blocks.

DMKFRERC - I want storage but do not extend the system to get it! Time driven sampling routines use this. They simply reschedule & try later.

DMKFRETR - returned storage goes directly to FRELS. Keeps subpools from growing

DMKFRETE - logoff cleanup. Branch to FRETR. Unextend if possible. Called by DMKUSP.

## Virtual Storage Management

Control Blocks:

VMBLOCK anchors segment table + SCBLOCK.

sectable points to page tables.

Pageable contains RFN of each page.

status

Abc points to swap table.

SWPTABLE - DASD location of slots → direct pointer if paged out. Indirect pointer points to SSBLOCK (swap set Block).

Hardware  
Defined

SCBLOK points to swap list + SSBLOKS for  
logical + phy swapping.  
SSBLOK - DASD location.

Demand Paging - two swap sets are  
prepage. The rest of the working set  
is built one page fault at a time. Needs  
a PIC 11, or 10.

CP can explicitly demand a page  
since it runs DAT off.

If CP ever pgm chks, it aband.

Pre-paging is initiated by the scheduler. Number  
of swap sets should be number of paths to  
swap sets for interactive users. For non-inter,  
prepage 1 swap set. The lowest vrit  
store addresses are brought in. No guarantee  
this is what is first need.

Page fault on a page in a swap  
set brings in the entire swap set.

Page status:

- Resident (see Page Table; bit 12 invalid bit)
  - Active  $\rightarrow$  bit 12 is 0
  - Reclaimable  $\rightarrow$  bit 12 is 1 + RPT is  $\neq \emptyset$
- In Transit  $\rightarrow$  flag in swap table
- Not Resident
  - Paged out - Page table entry  $\emptyset\emptyset\emptyset 8$
  - Never used -
  - Part of a swapped out set

## SEGTABLE - CR1

has flag byte in last byte of adr. This is ok since a page table must begin on a DW boundary (last 3 bits 0).

x'40' - seg ENQ

10 - seg mig

04 - seg prot

02 - seg common (MVS)

01 - segment invalid. PIC 10

### Examples:

F0123400 - in queue  
                    
          adr

F0123401 - invalid bit on; user in queue.

00000001 - never used. Only invalid bit is on.

Page Table Header - 16 bytes ahead of table. Has timestamp to see how often the segment is referenced. This is used for page migration.

For shared seg, has a count of # users. Also, pointer to share table (SHRTABLE). (-8)

Has pointer to swap table. (-4)

Entries are half words. Bits 0-11 are RFN. Bit 12 is invalid bit.

Bits 13 + 14 are extended addressing.

Bit 15 is used by page migrator.

This means migrator has been here if 08.

Swaptable - 16 byte header. Points to VMBLOK + pagtable.

Entries:

Flags → has double ref + chg bits to support 2K keys.

In-transit bit.

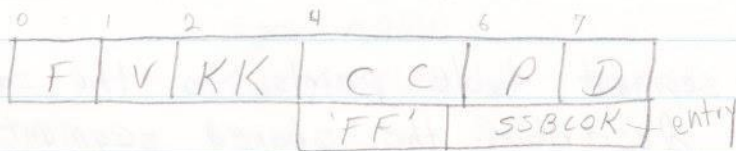
Shared sys

Permanently assigned. Migrator cannot move this. The cp nuc set this for himself + leaves NUC on best paging device.

virt page # (0-F)

Key - 2 bytes (for 2K pages!)

DW { CC - cyl #; If FF then pointer to SSBLOK  
 P - formatted 4K page slot  
 D - dev code - position in sysown list.  
 Begins with 0.



SC BLOK

Next Pointer } if on a list  
 Previous ↑

owners VMBLOK

# SSBLOKs }  
 ↑ 1st SSBLOK } if phy swapped  
 ↑ last "

↑ CORE table entry of first page on swap list (Logically, swap)

↑ last

↑ first trim

↑ last "

# pages in WS  
 # Trim

see VSM 110

## SSBLOK

Flag byte for in transit

Has SSBLOK entries to link virtual page number to each page slot (since they are not necessarily contiguous in the swap set).

## CP Virtual

CP does use seg & page tables - but manages it itself. Up to pageable Nuc,  $\text{virt adr} = \text{real adr}$ . Everything in CP has a constant virt stor address. It treats itself as a 16 meg virt. machine. Spool file buffer, etc. are created at init & then paged out.

3-2-88  
Hollis Johnson

## Shared Segment Management

Each user's segment table points to the same page tables describing the shared segment. The SHRTABLE lists the page table locations to be shared.

Shared segments are defined in DMKSNT.  
Attributes:

shared / Nonshared

↳ protected / unprotected

## Modules

- DMKCFG - Entered via logon, IPC CMS, or Diag 64. Calls DMKCFE.
- DMKVMA - ensure seg protection when VMA seg prot not avail.
- DMKATS - "unshares" changed pages + notifies users.

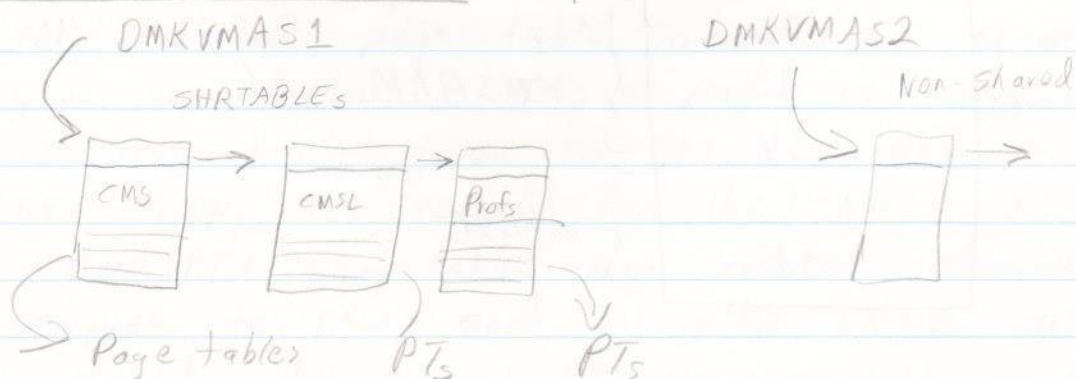
## Control Blocks

- VMABLOK - 1 / saved sys / user - chained off VMABLOK (usually multiple).
- SHRTABLE - 1 / shared sys or DCSS

VMABLOK - 16 bytes: ↑ next VMABLOK  
↑ SHRTABLE  
Sys Name  
↑ Non-shared sys

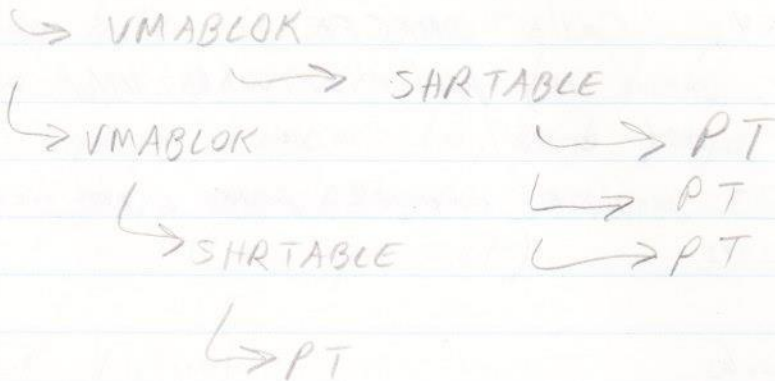
SHRTABLE: - Flag byte - indicates protection.  
sys name  
↑↓ pointers to shared + unshared tables on chain  
size in DWS  
# shared segs  
shared seg numbers  
Pointer to page table for each shared seg listed above.

Pointers in DMKVMA: System anchors -

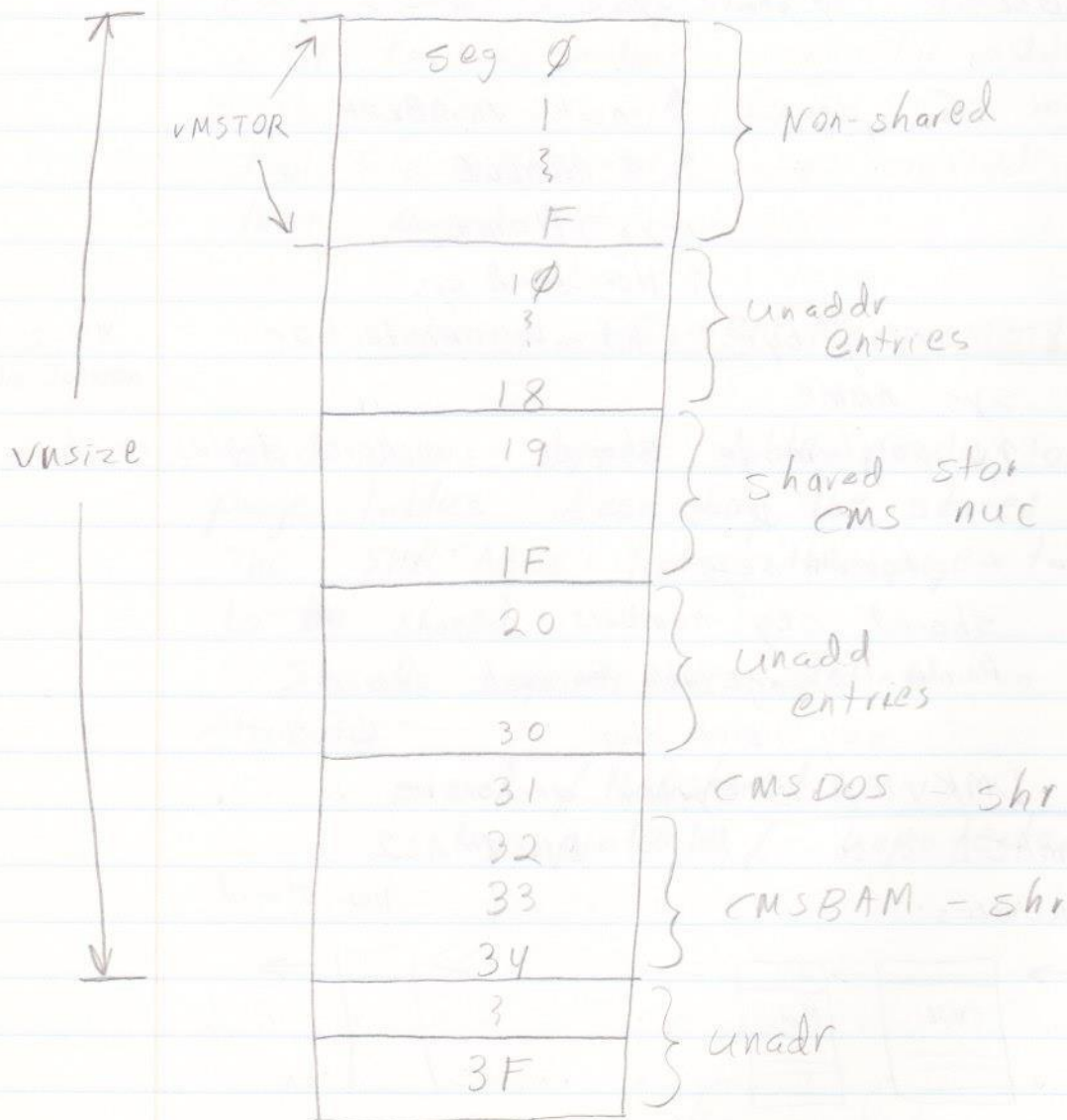


## Virtual Anchors

VMBLOK



### Segtable



VMABLOCK is not required when VMA seg prot is avail. PT is found through DMKVMA chains.

SHRTABLE is built when 1st user requests system + storage is freed when last user logs off.

DMKCFE - does setup + clean. Will build a SHRTABLE if required. Updates PSW + regs for user to begin exec.

DMKCFG - reset user on IPL cmd.

call DMKVMI if IPL by adr.

If IPL by name, access DMKSNT construct tables

DMKCFE to prep user.

DMKCFE - get PT + swap tables, updates them

- reads DMKSNT

- Determines shr/non-shr + anchors in VMABLOCK

- Builds new shrtable if req.

- modify PSW + Regs

- exit to disp.

w/o VMA protection, CP must check every shared page table whenever a user gives up control! If a page is chg, it is unshared + given to the user and a new copy is created for the shared users.

New PT + swap tables are created, changed entries copied to new, and all other entries are invalid. The original PT entry (corrupted) will be invalidated + the corresponding swaptable says



to read it in from DASD. Thus a new copy will be brought in whenever a user tries to reference the page. The new SWPTABLE entry is invalid, indicating a copy of the page does not yet exist on Aux. As the changing user runs, the other entries (now blank) in the page table will be built.

AP/MP - each CPU has a set of page + swap tables.

Hardware Seg Prot - bit 29 in segment table entry (X'04). Causes seg prot if change attempted. PSA flag indicates if seg prot is active. VMABLOCKS are not used.

Sharing reduces paging because it reduces the number of pages trying to reside in memory.

Hollis

## Paging / Swapping

At Q-Drop the PT entries are invalidated + separated into Trim + working set. These are anchored in an SCBLOCK. Trim is placed on Flushlist (DMKPTRU1). WS is on swaplist. PTEs of both will be in the form XXX8. Pages are invalid, but still in + therefore RPN is ok. Coretable entries will not have a VMBLOCK pointer; instead it will have pointers to flush or swap list, plus flags.

When user comes in-queue & references a page w/ invalid bit a PIC 11 occurs.

DMKPRGIN detects PIC 11 + calls DMKPTRAN which returns CC 0 or 2. 0 = page is now avail. CC 2 = means virt adr was out of range & causes PIC 5, addressing exception.

To reclaim, turn off bit 12 in PTE, update coretable w/ VMBLOCK, + repair chain to remove PTE.

DMKPTRAN exec CRA - load real address.

CC 0 = pg in storage → then exit

1 = seg. invalid →

2 = PTE invalid

3 = addr exception

This test governs what PRGIN needs to do in order to make the page available.

Two types of paging: 1) Demand (PIC 11) which is either page or swap fault. 2) Prepaging. CP doesn't PIC 11. It uses the TRAN macro which causes DMKPTRAN to run.

Pages will be either on flush, swap, or free list, or else out on aux.

Free list pages don't need to be paged out.

If pages on swap list (log swap) all will be made available when user moves in-queue.

Page Frame Status: (on PTE)

Active - XXX2

Reclaimable - XXX8 (+ coretable tells which list)

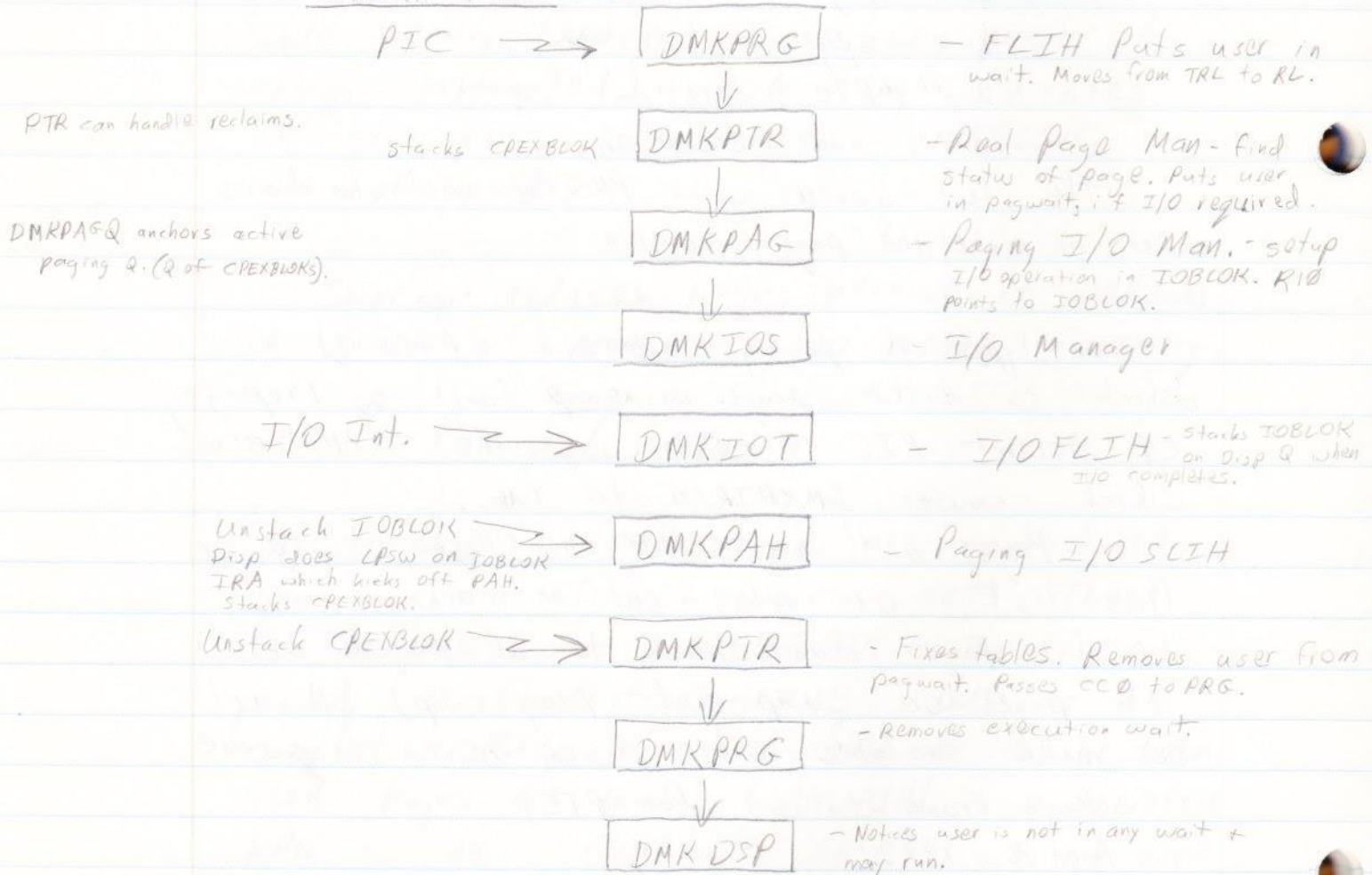
In Transit is a special case. coretable flag has CCW op code. PTE = XXX8  
 Swap table shows intransit.

If frame is paged out, PTE = 0008.  
 If never ref. PTE = 0008 + swap table is all zeros.

### CP Explicit

R2 has parameters + DMKPTRAN is called. On return, real address will be in R2.

#### Overall Flow



DMKPTR does a local real addr. (LRA).

This causes addr trans to occur & the real addr is returned in a reg if available & CC =  $\emptyset$ .

CC = 2 = PTE invalid - most common.  $\rightarrow$  Reclaimable, In Transit, Out.

1 = STE " - then go to DMKSTR (seg Tab Rebuild)

3 = Virt address is out of range.

LRA R2,  $\emptyset$ (R1) lets CP to do a translation without DAT being on.

DMKPRGIN - entered via PSW swap. If not a PTC 10 or 11 goto DMKPRV or FPS for simulation. Put user in Exec wait. SVC to DMKPTRAN. If user is EC mode (VMBLOCK virt PSW) then goto DMKVAT for use of shadow tables.

DMKPTRAN - Bring page in. This is invoked with options (Bring, Defer, VFault usually).

R11 = VMBLOCK.

Reclaimable: (cc=2 from LRA)

If page is on freelist:

increment reclaim count (DMKPTRPR + ER)

decrement # free pages (DMKPTRFN or N2)

set storage keys

Reset ref & chg bit (leave chg on if on)

Update coretable entry:

VMBLOCK in 1st word

Lock Ct in 2nd "

Set flag

Leave PTE & SWPTE alone in CTE

Put real adr in R2

Return to caller w/ CC =  $\emptyset$

common even  
if page is  
on flushlist

In Transit : (cc=2 from LRA)

check swaptable + see page is intransit.

Defer is <sup>not</sup> specified : (only if CP invokes TRAN + can't wait)

Build CPEXBLOK

Exit cc=1

Defer specified : (CP is not demanding the page immediately) This is usual

Add pag wt ct

Place user in page wait

Build CPEX + chain to existing

CPEXBLOK because this page is already on the way in.

Exit to DMKPAG

Not Resident : (cc=2 from LRA)

Add pag wait count

Place user in pagewait since I/O req.

acquire free frame from freelist

update CTE

Flag swap table in transit

If first ref, then return to caller cc=0

Else, page-in necessary.

stack CPEXBLOK

call DMKSEC

PTRAN Summary

LRA → 0 - resident → lock if requested, + exit cc=0

1 - seg. Table Rebuild (STR) + try LRA again

2 - Page Table Exception

- Reclaimable

- In Transit

- Not Resident

3 - Exit cc=2

your R3 will then get control when page is back. It is the IRA for the CPEXBLOK.

DMKPTR can also lock (DMKPTRCK) or unlock (UL) a page, or request an available page from freelist (FR). - this last emulates a first-ret page fault.

## Swapping

Swap In - data is in a swap set on aux.

DMKPTRAN still handles + determines if paged out or part of a swap set. In swpTable, field SWPCYL will be

FF1↑SSBLOK
('FF' in first byte is flag)

if swapped +

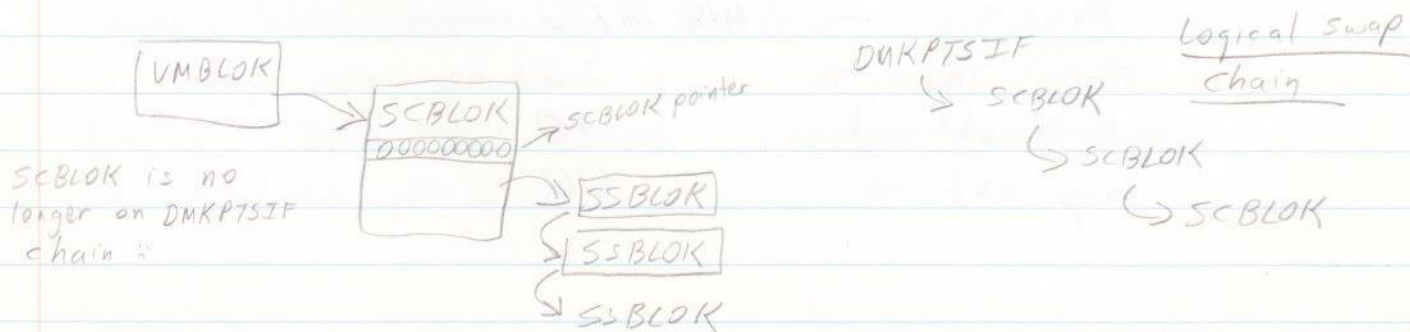
cccc pp dd

if paged.

Logical Swap-In - DMKSWA does pre-page or logical swap in. Check SCBLOK to see which.

SCBLOK is no longer on chain when physically swapped. SCBLOK had been chained with other SCBLOK at Q-Drop. DMKPTSIF is anchor for these. Also, SCBLOK now points to SSBLOKs. ↑↓ SCBLOKs are zeroed when removed from chain.

Then remove CTE from swap list chain, put VMBLOK↑ back, reset CTE, + hardware keys.



Pre-Paging: do logical swap in of any remaining pages. Bring in # of swap sets specified using lowest addresses of sets actually referenced during last execution.

Swap Out - swap out occurs when DMKSELECT need pages for freelist:

- 1) Garbage
- 2) Old userswap
- 3) Flush
- 4) New users
- 5) Coretable scan

Get SCBLOK from non-inter queue. If interact user still logically swapped for DMKSWATIB (20sec), then move him to non-inter queue.

Call DMKFREE for storage to build SSBLOK. Chain SSBLOK to SCBLOK.

call DMKPGT to get slot on DASD.

Call DMKPSW to get rid of old slots.

Get CPEX to do swap out.

If more pages must move, build more SSBLOKS.

CPEXBLOKS are chained off:

Page in	-	DMKPTRRQ
Swap in	-	IQ
Page out	-	WQ
Swap out	-	OQ

DMKSWA gets control from CPEXBLOK.

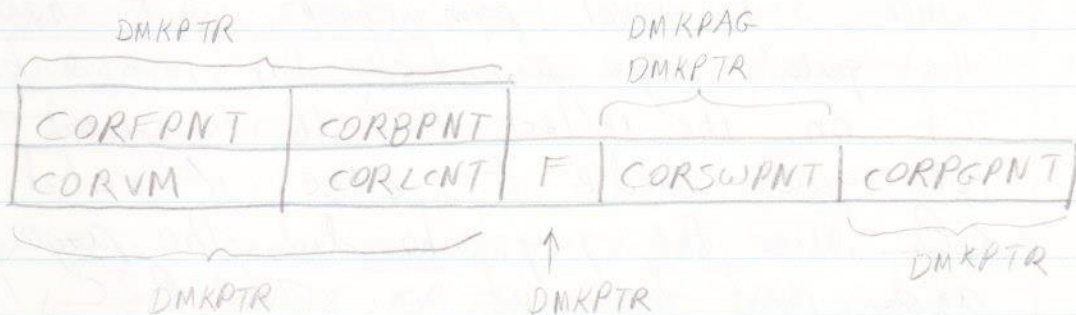
Set swapout flag in SSBLOK.

call DMKPTS to add page to freelist queue.

### Real I/O

DMKPAG calls DMKIOS which queues IOBLOK on RDEVBLOK. All IOBLOKs are queued (+ sorted) on RDEVTIOB. If multiple IOBLOKs required same cyl, then only one seek is required + will be executed together. The active IOBLOK is anchored @ RDEVAIOB.

### Coretable Entry



3-3-88  
Jim Turner

### Shadow Tables

Shadow Tables are simply another set of segment and page tables. These are used for second-level V=V virtual storage machines. They can eliminate double paging. Shadow Tables map guest virt to CP real. The STOR gets loaded into CRL. Shadow Tables involve high overhead. CP SP can map six shadow tables per user (HPO handles 15/user). Each shadow table maps an address space.





For non-paging guest, CP does everything. The guest never receives PIC 11. Its believes memory is real.

Shadow Tables have a combination of into ("real" RFN) to map guest virtual to CP real.

Guest cannot modify segment + page tables w/o issuing privileged inst. (Load CR 1, SIOF, IPTÉ, (for paging) which is prohibited + CP intercepts the pgm chk + knows to check. The shadow tables are checked prior to the user being dispatched (via load CR1).

Shadow Table Creation:

Step 1: Copy guest Seg + Pag tables to shadow tables.

Step 2: Check each entry + do addr. trans to see where CP really put the page. Imagine the overhead.

<u>Guest Virt</u>		<u>CP Virt</u>	<u>Guest Real</u>	<u>CP Real</u>	<u>Shadow</u>
020000	} Guest Does This Trans	}	12000	456000	456000
021000			23000	N/A	N/A
022000			N/A		
023000			3400	789000	789000

CP will always build a shadow table for a EC, DAT V=V machine even if it has to give up one already built! and create a new one just to dispatch the user!

DMKVAT determines who did the page out on a PIC 11 + will go to either DMKPTR or PRG.

VMBLOCK → ECBLOCK (this should look like an adr).

Control status flag:

shadow tables present? valid?

Has CR 0 changed (virtual CR 0 means guest has changed architecture).

Is CR 0 valid? (NO, pgm chk user) <sup>IF BC mode</sup> CR 0 not chk.

Is user in EC mode?

### ECBLOCK

save areas for virt CRs.

" " " shadow control Regs 0 + 1

(These point to the tables running this address space. These will be loaded in real CRs).

↑ first STOBLOCK (Default 3, SET STMODE x)  
count of STOBLOCKS.

### STOBLOCK

↑ pointers to next + last STOBLOCK  
(last is in case one needs to be trashed).

Virt + Shadow CR 1

size of VV = VR area (guest nuc)

last VV = VR page table.

Shadow Segment Table (not pointer).

↳ points to shadow page tables.

Common Shadow Page Tables - various tables under a virt machine can share common page tables. - This is for the commonly addressed portions of MVS address space. The limits must be exactly correct in the CP command or MVS address spaces will corrupt MVS storage!

SET STBYPASS  
tricky to use or  
worth it.

DMKVAT - Virt Adv Trans.

Does table maintenance -

Builds tables if bit 5 + 12 on.

Checks when CR0 or 4 is reloaded.

Any user page operation.

Guest PTLB, or IPTE instructions.

Allocates + initializes as soon as guest enters EC DAT mode.

When the user leaves EC mode, the storage is freed (ok to leave DAT mode).

VAT determines who (guest/CP) did a page out when a PIC 10 or 11 occurs. Entered from DMKPRG if bits 12 + 5 on. Call DMKPTR if CP invalidated entry. Also, give user the PIC.

DMKVAU - utility routine.

Does the translations of the guest tables to build the shadow tables. (3rd to 2nd level and 3rd to 1st).

Pseudo Page Faults - part of DOS + VSI hand-shaking (SET PAGEX ON). Allows guest to suffer multiple page faults. If on, the user is not placed in page wait until the page arrives. The guest can run other work (VM= YES in guest gen). CP passes a PIC 14 (not in architecture!) + the guest dispatches other work. A second PIC 14 means the page is now in. PG-BLOCK ← ECBLOCK holds the virtual address of the page to be brought in.

## DASD Storage Management

CP tries to choose the best spot on disk to place a page.

Spooling is really non-preferred paging as far as management is concerned.

DMKPGT - Alloc DASD space (slot)

PGU - Dealloc " "

TDK - Alloc/Dealloc TDISK (using some control blocks)

VDG - @ init builds Allocbloks + Recbloks.

OWNDLIST - List of DW for each vol in the sysown list holding volser + pointer into RDEVBLK. Resides in DMKSYS. Order of sysown list determines the order of volume use (HPO overrides with sys page)!

RDEVBLK - describes device + if CP owned points to Allocblok. RDEVBLK are anchored in PSA at ARIODV.

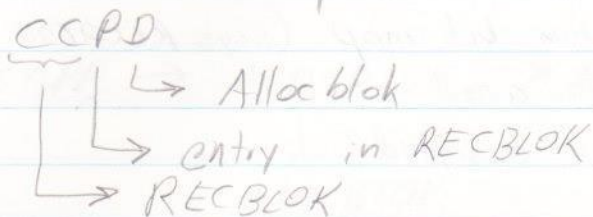
Allocblok - three types:

Full Volumes - describes CP space within pack. Used only to Dealloc slots. ↑ from RDEV.

Area Allocblok - describes use of CP area + is used to allocate. ↑ from sysplist

TDisk - describes TDISK space + use.

RECBLOK - represents a cylinder. In CCPD, CC points to RECBLOK & P points to entry therein. CP checks each RECBLOK when looking for a free slot.



SYSPLIST - created from syspag macro.

Defines use of each volume. Defines hierarchy for paging. Swapping → pref pag → gen pag → pag mig → Temp.

Be sure to round-robin among devices in same class.

Primary point of control is SYSPAG + SYSOWN macros. 4% of CPU cycles should be for paging. If more, then the dispatcher is effected. Take care here. IND command lists paging as a % of cycles; also SMART + VM MAP.

ALOCBLOK -

type flag (SW, PP, PG) may be combined for full-pack Allocblok.

↑ next deallocation Allocblok (PGU follows this)

Cyl count in use / in area (how many empty cylinders.

cyts in use may have only one page. A used cyl has a RECBLOK).

# pages alloc

# pages in RECBLOKs

↑ RDEV Blok

anchors for paging + spooling, RECBLOK chains

↑ next allocation ALOCBLOK

cyl pointers

↑ to current RECBLOK (where the arm is).

lots of count fields

cyl allocation bit map (says RECBLOK exists)

0 = avail

1 = assigned

RECBLOK - ↑ from ALOCBLOK. Only exist if needed.

↑↓ Pointers

cyl #

Pages in use

Pages in cyl

Page alloc bit map

0 = avail

1 = assigned

→ Flag - type of device + use

SYSPLIST - created from SYSPAG macro.

DMKSYSSW  
PP  
PG

Flag

↑ next lower level in hierarchy

↑ previous sysplist

↑ allocblok currently being used

count fields

↑ sysplist extension

Then for each volume:

volser

Dev number

starting cyl

ending cyl

area type

Pointers: Allocation chain

DMKSYSW → ALOCBLOK

↳ RECBLOK → swap space

↓  
RECBLOK → swap space

DMKSYSPP → ALOCBLOK → RECBLOK → Page space

↓  
ALOCBLOK

↓  
RECBLOK → Page space

DMKSYSPS → ALOCBLOK → RECBLOK → spool space

↓  
RECBLOK → spool space

Deallocation Chain:

RDEVBLOK → ALOCBLOK

↳ RECBLOK → DASD Space

Area Usage:

<u>FMT</u>	<u>Syspag</u>	<u>Anchor</u>	<u>Type</u>
Page	SW	DMKSYSW	- swap
Page	PP	PP	- pref paging
Page	PG	PG	- gen paging
Page	PM	PM	- migration
Temp	PS	PS	- spool
Dump	DU	DP	- Dump

At CP Init:

DMKALO (from DMKCP) builds full pack ALOCBLOK.

DMKVDG creates area ALOCBLOKs.

Chains in hierarchy order. RECBLOKs are pre-allocated for all space but spool + migration.



DMKPST creates ALOCBLOKs + RECBLOKs for paging storage + chains to sysPList.

### Slot selection

Based on -

- level
  - volume
  - slot
1. syspag hierarchy (type determines where to try first).
  2. N-Select - what vol. on level
  3. Moving cursor - where on volume.

Start at lowest cyl + move to highest vol., then reverse + sweep back. This favors interactive users.

Modified Moving Cursor - beginning to end + back to beginning. Used for spool, migration + dump.

N-select - determines how long to use a single ALOCBLOK. When count (ALOCCTMI) reaches 0, go on to next ALOCBLOK.

<u>Dev + Type</u>	<u>ALOCCTMI</u>
SW, PM, PS, DU	1
PP, PG	1
3380	10
3350	8

### Page Out Blocking (HPO Only)

HPO tries to conserve I/O operations. Multiple IOBLOKs are stacked on an RDEVBLOK + CP can group some together. When a page is ready to go out, CP will queue it on a dummy RDEVBLOK until reaching N-Select number. Then they are turned over to the I/O manager. Now a better chance of

anchored in  
DMKPAG

chaining these together. <sup>+ going on the same cyl.</sup> This process is not done for cached devices.

## TDISK Space

DMKTDK does Alloc/dealloc. It doesn't use RECBLOCKS. Called from DMKVDA and VDR. In SP, tdisk is by device addr in ascending order. In HPO, it is in the sysown list bottom to top. To get better use, have CP search the smallest space volume first.

That way, small requests for tdisk won't fragment the large spaces.

No hierarchy of devices. The user request is specific.

DMKZTD clears deallocated space if CLEAR=YES is specified in DMKSYS. CMSBAT machines should have mini-disks carved out of TDISK.

TDISKS - ALOCBLOCKs are anchored @

DMKPGTAS - 2305

8A - 3380

A0 - 3330

50 - 3350

A4 - 2314

Jim

Page + Suptable Migration

Page Mig - purpose: to free pref. space

Suptable Mig - purpose: to recover high free.

Once suptables are out, pseudo page + suptables are created in order to page in the suptables.

Only 184 bytes of high free are needed to map a 1 64K segment.

Page Migration uses regular paging routines: simply page in from perf + page out to non-perf. (Same for swap).

DMKSTP - performance monitor, checks flags + when discovers migration flag on, it stacks CPXBLK to run DMKPGM to mig a page or DMKSWM to mig a swap slot. Threshold for mig is set by SRM PGFULL. When exceeded, user is selected for migration, along w/ # of pages.

To select: loop through user's seg + pag table. Mig if seg tab invalid, not part of shared seg, + not part of swap set.

Mig runs constantly in HPO 4.2 + mig a little at a time.

Save Mig #140

### Swptable Migration :

Criteria : Paging > 4% of cycles,

Must be short of high free (Extends taking place).

User not in runlist.

DMKDSP unstacks CPEXBLOCK & calls DMKSTR which move swptables to page buffer.

At PIC 10, DMKPTR call DMKSTR which checks to see if the swptable is out. If so, bring it back.

Swptable migration is indicated in the segment table entry:

F0000011

F0000041

F0000019 - Priet space Mig

F0000013 - Gen space Mig

Pseudo pagetable has entries for each segment in sequential order.

VMBLOCK, VMSWPMIG → Pseudo page table (in invalid bit)

↳ pseudo swptable

has 16 entries

seg 0	seg 1
seg 2	seg 3
seg 4	...

0	CCPD
1	CCPD

3-4-88

## Timer Management

Jim

TOD Clock - come with 370. Incrementing DW register. bits 52-63 are not used. Bit 51 increment 1  $\mu$ sec. Bit 31 increments 1.0486 sec. Set by SCK instruction. Obtained via STCK. Always runs.

Clock Comparator - DW "system alarm clock". Same format as TOD. Set by SCKC + read by STCKC. Does not increment. Continuously compared of TOD + when  $<$  causes Ext Int 1004. Int are not kept pending. Be sure to enable for Ext Int before resetting. Also, be sure to reset while handling an Ext 1004.

CPU Timer - same format as TOD. Used as a stopwatch. Positive value is loaded + decrements at bit 51. Int is generated when it goes negative (Ext 1005). Set by SPT, read by STPT. Int are not kept pending. Does Not decrement when CPU is stopped.

Interval Timer - 360 Arch. Fullword in PSA. Bit 23 is decremented 300/sec. Not as accurate as CPU Timer. Stops when CPU is stopped. This decrements + generates Ext Int 0080 when negative. Int are kept pending.

TOD is based on 1 Jan 1900 GMT.  
Difference between local + GMT is in DMKSYSTZ.  
The previous local midnight is stored at PSA TODATE.  
Latest STCK is kept DMKSYSCK. CP uses the  
TOD to schedule deadline priority calculations +  
controls placement in the runlist. Uses bits 7-38.

Clock Comp. is used to schedule delayed execution.  
DMKSTP, the performance monitor uses this. The TRQBLOCKS  
are chained in order off DMKSCHTQ (in the  
scheduler). Top request is current. User requests  
are also stacked here.  $\rightarrow$  the one currently  
loaded in the register. At pop, the TRQ gets  
moved to another queue for execution. The chain  
will never be empty. The SCHTQ is itself a  
TRQ with a TOD value of 19 Sept. 2043.

CPU Timer is used to control access to the  
runlist.

Int. Timer is used only for Dispatcher timeslice  
+ is used to control access to the processor.

Timeslice is based on this + since it decrements  
300/sec., the timeslice will not be shorter than  
 $1/300$  sec.

SET DSPSLICE  $nnnn$  controls the dispatcher  
timeslice. Lower this to help interactive users. Lower  
as long as TVRation is not effected.

PSA values: Pagewait, IOWait, Fdllewait.  
Initially set to 7F—F. CP loads this when  
waiting, enables for interrupts, + when int comes  
stores value back in PSA. CPSTAT3 tells  
which field is in CPU Timer.

Determining whether I/O or Pag wait is real hooie. If more than 4% of CPU is for paging, then wait must have been for paging!

### PSA

Real Interval Timer  $\times 50$

TOD as MM/DD/YY

last midnight TOD DW

Idle

Pag

I/O

Problem State

CP Running Status

} Total System

### VMBLOCKs

ECBLOCK  $\uparrow$  (permits use of clock comp + CPU Timer)

Running status flags

Pending interrupt flag

Interval Timer Value

Virt CPU Time used - Problem State Time

Time left in queue

Total time in supervisor - Is really Total Time

$\uparrow$  TRQBLOCK holding real CPU Timer Value in virt CPU Timer is in real.

CR6 savearea

$\uparrow$  chain of  $\times$ INTBLOCKs (stored while user is disc for EXT INT)

$\uparrow$  TRQ for delayed sleep or lagoff

$\uparrow$  TRQ for Queue Drop Delay Timer

## ECBLOK

Virt CPU Time

↑ TRQ for CPU Time } Only one ↑ for each  
↑ TRQ for Clock Comp } because that is all  
that is architecturally  
available.

## IRQBLOK

Clock Comp Value for Int.

TOD that TRQ was queued

↑ VMBLOK

IRA - must be non-pageable.

Real CPU Time Value if the Virt CPU

Timer is loaded in real CPU timer

## XINTBLOK

# of Int

Ext Int code

CR 0 mask (bits 16-31)

Ext Int parameter word.

DMKTMR manages timing functions - simulates timing + interrupts. To simulate inst, stack a TRQ - see if user already has one, yes then change it; no, then stack one.

To simulate int, go into VMBLOK + add XINTBLOK.

Also calculates users' V + T Time.



DMKPRG receives int. Try Fast Path Simulation (DMKFPS). If not able, go to DMKTMR.

### Inst. Simulation -

If user not EC mode, pgm chk.

If address not D<sub>w</sub> aligned, spec exception.

Bad addr - addr. exception

Check storage key

Simulate

DMKDSP

SET TOD Clock - returns CC-0 + does nothing.

### Set CPU Timer -

DMKDSP updates VMVTIME + VMTT<sub>IME</sub>, Virtual Interval Timer, reflect interval time int. to user, + acct. for system wait.

Virtual Int Timer - is in page 0 of the user's virt. storage. SET Timer ON or SET Timer Real causes updating.

Int Timer does not run if machine is stopped. When user not running virt. machine is stopped.

Int Timer is decremented with value of timeslice (less any CP READ time). Update occurs prior to dispatch. SET TIMER REAL will cause non-running time to be decremented (MORwatch!).

With VMA, user's page 0 gets locked in storage + ucode handles decrementing.

With SET TIMER ON you cannot set a Int Timer Value & then go into an enabled wait. You won't be dispatched until an int comes in & the interval timer field won't be update until you are dispatched & thus will never issue an interrupt.

Scheduler maintains TRQ chain. To stack a TRQ, put address of TRQ in R10 & call routine. Scheduler does no validity checking of TRQ. It will stack an invalid TRQ & crash the system.

New TRQs are inserted into the chain in chronologic order. If the same TRQ is stacked twice, scheduler will loop trying to insert it in the chain.

Only one chain even on a MP. (ACTIVEQ)

### Pseudo Timers

A device may be defined as a timer and a SIO will bring back time values.

Also, DIAG PC brings back DW CPU timer.

Both give MM/DD/YY & HH:MM:SEC & V & T CPU Time.

### DIAG 70

Defines field for virt Time & current TOD value. CP locks the page & these fields are updated prior to dispatch. No interrupt is required to access the values since they live in user virtual.

## Spooling

Lots of control blocks + modules.

DMKSPL - spool file manager

DMKRSP - real spool manager

Finding Devices:

PSA

↑ Punch Table

↑ Prt "

↑ RDR Table

All point to DMKRIO for # entries,  
Displacement into RDEVBLK chain, + CUG.

Spool File attr.

class - determines device + priority

name - 12 char

type - 12 char

spool id - 1 - 9900

Rec Count - 1 - 16,000,000 (cc is a rec)

Copy Count - 1 - 255

Real Stor Man.

Allocates frame at SIO.

Locked during data transfer. - only the  
frames in the current SIO are locked.

CP can only have a max of 1280 spool  
file buffers. Thus only 1280 spool operations  
can be active at once. Else, P&G&P&S.

All spool I/O is handled by DMKPAG and uses separate RECBLOK chains exclusively for spooling.

SFBLOK - Spool File Block

VSPCTL - virt sp CB - only around during spooling

VSPXBLOK - ext. to VDEVBLOK

VPRXBLOK - virt prt extension

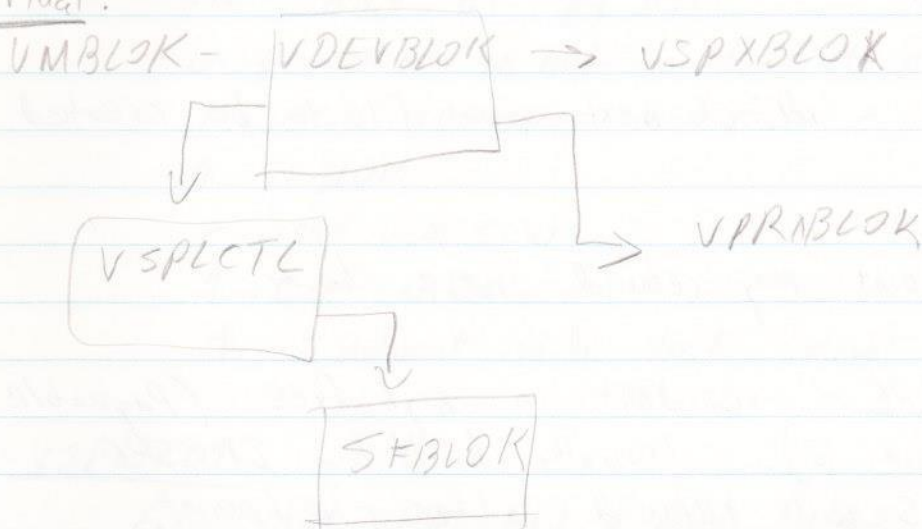
RSPCTL - Real sp CB

RSPXBLOK - ... extension - 1/real prt or Pun

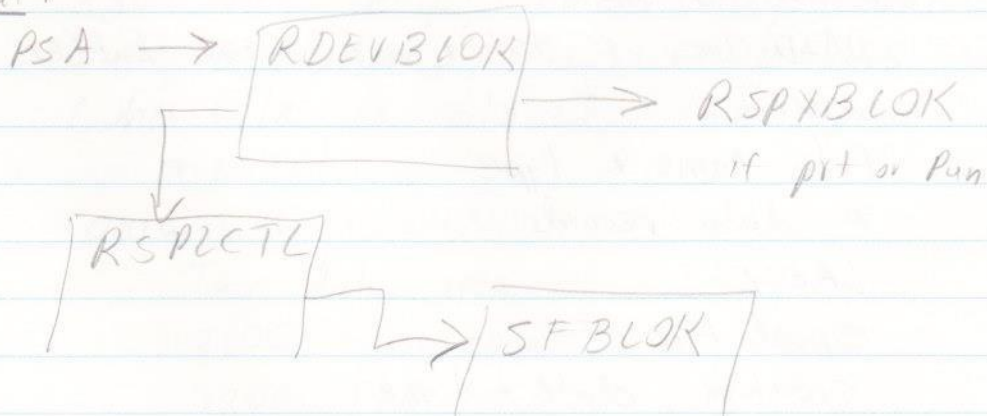
SHQBLOK - Sp hold Q Block

SPLINK - Spool Buffers.

Virtual:



Real:



DMKRSP - chains of RDR, prt, & pun files. Has execatchers. Also has a chain of files to be deleted. Chains point to SFBLKs. PSA points to these anchors in DMRSP.

### DMKRSP

PRT

PUN

RDR

DEL - to be deleted

ACNT -

DUMP - dump spool files

HOLD -

SPT -

ID - id of next spool file to be created.

:

Checkpoint may rebuild these chains.

SFBLK - resident in high free (Pageable in release 5). Usually built by DMKSPC. DMKCKS will rebuild after checkpoint.

↑ next SFBLK

DASD loc. of 1st + last page buffer. CCPD  
userid

file name & type

# data records

LRECL

Spool File #

creation date & time

# copies

spo file class

Distribution code

chkpt map # - what slot on what cyl.

form #

Control flags

### VSPLECTC

Exist for open files (temporary)

Used to chain SFBLKs to VDEVBLK.

Built by DMKVSP (or DRD for dumps).

Contains info need to process file:

virt addr of user's ccw

QASD loc of current page buffer

virt add of pg buf

Displacement in buf to start of next rec

# recs remaining in buf

↑ SFBLK

current user ccw

↑ Buf area

↑ Indirect adr list (might cross pg boundary)

### VSPYBLCK

Exists for open prt + pun files;

3800 prt setup; holds tag info.

Built by DMKCST (tag) + DMKVDS (3800).

When file is closed, the specifications are set.

Contains:

Tag data area

3800 flash count

3800 FCB

3800 char arrangement Table

3800 copy mod

Name of flash overlay.

RSPICTEL - exists for open real files.  
Restart info (pit runs out of paper).

Contains:

Restart CAW - CCW adr

DAJD loc of cur pg

Real Adr of pag but

virt Adr " "

↓ SFBLOCK

RSPIBLOCK - ext to RDEVIBLOCK. for every  
real Pit + Pun. Has form info, etc.

SHQBLOCK - exist for each user w/ files  
on HOLD. Permits future files to  
be held for a userid.

SPLINK - this is the spool data:

header + (optional) tail.

Chained from SFBLOCK.

Active SPLINK lives in DPA +  
while there chained from RSPICTEL +  
VSPICTEL.

Built by DMKSPL.

Contains CCW to drive real device  
and data.

Only first SPLINK has Tag, FCB, +  
trailer w/ date/time stamp +  
3800 info.

↑ + ↓ - SPLINK

Note: cannot backspace a 3800 printer.

Data area of SPLINK is 4080 bytes + is not compressed. Data is actually interspersed with the CCWs. First CCW is NO-op with bit 35 (ship) = 0. RSCS looks for this + sees that 136 bytes of TAG exist. Next is TIC CCW + is FCB CCW, + then write, tic, data ... At end of CCW chain is another NO-op so CE/DE can come back simultaneously. (This is 370 arch).

The CCW addresses are virtual. Write CCW will be 00000C. Tic will be C + length of data.

### Spool File Recovery

Cold - nothing recovered

Warm - SFBLOCKS

RECBLCK

SHQBLOCKS

} all recovered

CP simply does chaining in light free.

CKPT - Read every rec of every spool file to build RECBLCKs. Does everything sequentially. No overlapping I/O.

If invalid chain found, recovery stops.

Open Acct records lost.

Force - if an invalid file is found it is thrown out + the chain is repaired.

No indication of which files lost.

Really no advantage to using CKPT except you know that no files were lost.



On shutdown, all info is taken out of high free to build RECBLOCKS + written to warm start area.

Chkpt - each entry is 128 bytes. (32 rec/page)  
1-6 Chkpt maps. Each entry in the map describes the respective chkpt 128 byte entry.

Table:  $2^{16}$  'FFFF' = empty

'AAAA' = SHQBLOCK

'EDIE' = status of del DIE

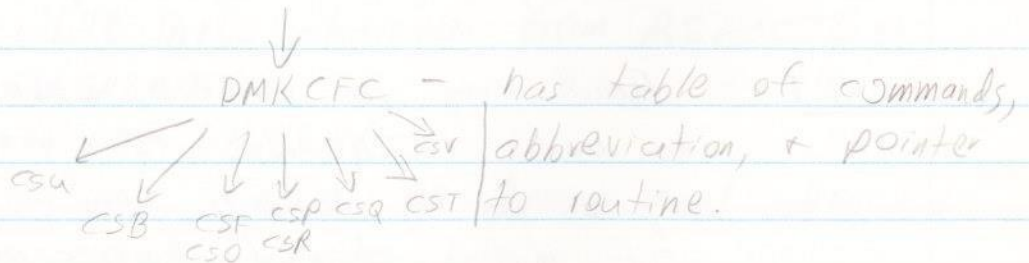
'FEFE' = End of Table

To delete, find map containing entry, read into storage, then read the record in, + invalidate chkpt entry.

MC

DMKRSPM1 - M6 points to virtual address of the chkpt maps. This is in CP's virtual storage.

console cmd → DMKCFM



DMKSPC - Build the SFBLOCKS + do the chaining. DMKSPK handles deletes. Two routines: 1 stacks SFBLOCK on del chain + another that deletes.

DMKRSPM1 + M6  
point to the starting  
ending cylinders.

DMKSPC - flow:

To open virtual spool file:

Build VSPCTL

Get a page (DMKPGU)

Get DASD Buffer (DMKPGT)

Get spool # (DMKCKT)

Lock page (DMKPTR)

Init SPLINK

IF 3800 get large VPRXBLOK

To close:

Update SFBLK

Read 1st SPLINK (DMKRPA)

chkpt (DMKCKS)

Validate + send msg.

DMKSPK - Delete routine:

To delete:

run chain of del SFBLKs.

Reclaim DASD space (DMKPGU or DRD) single-threaded

Delete SFBLK

Give back spool'd.

User is placed in I/O wait while spool file management takes place - even if no physical I/O occurs.

DMKRSP - Real Spool Manager

Holds spool anchors + counters:

↑ RDR, PRT, + PUN Chains

↑ Acct card Buffers

↑ Delete chain

↑ Spool Tape Delete Chain

Spool ID counter

↑ SHQBLOK

↑ active monitor SFBLOK

Virt addr of Chkpoint Maps

:

RSP reads the SPLINK in + adjusts the write CCWs. It knows where it put the SPLINK + thus changes the CCW addresses. It maintains an address counter as it works the way down the SPLINK. Then it alters the CAW to point 'x'10' into the SPLINK.

3-7-88

## Nucleus Creation

Gary Weinhdd

414-347-2390

VMFLOAD Loadlist CTRLfile - exec the builds a nuc in the reader. The loadlist is a list of modules. The CTRLfile determines which version of various modules to INCLUDE. Output is an IPLed card-image file. Movefile to tape from RDR.

Loadlist determines the order of modules in the nucleus.

## Building CP

DMKLDDEE - DMKLDDEE LOADER - stand alone loader  
for deck (0/4 bootstrap cards)

DMKPSA - first nuc mod

DMKSEC - if V=R present (1 card: sec adr of  
where to set the location counter as we  
build this nucleus).

DMKACR = order of resident nuc modules does not  
really matter, except DMKCPĒ which must  
be the last.

DMK CCH -

CPE - end of nuc

TAP - 1st pageable mod. Aligned on page boundary. Must be first.

VMD  
⋮ } Pageable modules

DMKSAV marks the upper limit of pageable  
modules that will be available before initialization  
of the paging subsystem.

CKP - now split:

CKD - } Don't move, change, or add to these.

CKN - }

LDT - Loader Terminator. Must be last card.

TPL DDC loads and executes the loader: 4 card bootstrap  
loads DMKLDDEE, reads CP text decks, links  
all modules (resolves addresses), creates an  
executable CP nuc in storage, + transfers control  
to DMKSAVNC. Any errors are printed. Also,  
entry points are printed. Comment cards are  
considered errors. The result is the NUCMAP.

The LDT card has an entry point which must  
have been resolved. This is DMKSAVNC.

Loader wait states are documented separately  
from standard WSCs.

DMKSAV NC writes the Nuc to disk (using rol & loc. found in DMKSYS). The current TOD is placed in DMKCP1. Ends in disa wait state '12'.

DMKSAV saves the data needed to IPL in order to create the IPL text. DMKCKP becomes record 2. CKP must be able to find the rest of the nucleus. The bootstrap gets DMKCKP running.

Nuc load termination:

12 - normal

10 - wrong volser

11 - fatal I/O error

0A - mch chk

14 - Storage failure

## Initialization

IPL text:

Rec 1, 24 bytes:

IPL PSW = 000c0000 00001000

CCWs = 00001000 00001000 Read

0 ——— 0 ——— 0

At end of chn pgm, the psw at loc 0 is executed. The Read CCW puts DMKCKP (1 page) into storage loc '1000'. The second CCW is not used. Chain bit is not on. PSW begins execution at DMKCKP.

DMKCKP checks PSA CPID to see if the value is CPCP or WARM to decide if warm start data should be saved. If either then it will try to save warm state data.

Warmstart may also be written at shutdown or abend. At abend, after dump is taken DMKDMPRS will do an I/O to bring in DMKCRP and start it running.

If CPCR then save warmstart data and end WSC = 8. If anything else, keep coming up. When coming up DMKSAVRS uses the warm start data.

PSA area checked by DMKSAV:

CPCR - save warmstart + shutdown

Warm - " " + restart

xxxx - normal IPL on cold machine

No warmstart data saved.

These fields are then changed to:

WRM - for restart

COLD - for normal IPL.

Init Sequence:

DMKCRP reads in CKD, CKF, CKH, CKM, + CKN right after CKP. DMKSAV gets loaded to its correct location in NUC. Throws away all I/O interrupts (for users pounding their keyboards). Then check to see if warm start.

Warmstart data: save -

- avail console

- addresses of en terms

- Acct info

- log msgs

- SFBLKs (from DMKRSP ↑)

- status of UR devices (from PSA ↑)

- RECBLOKs (from SYSPCIST + ALOCBLK)

- Spool Hold Queue Blocks. (SHQBLOKs)

An IPL over a running system w/o clear option will save warmstart + shutdown.

DMKSAV - loads nuc up to itself but no further. Throws out all I/O ints except from sysres. Exits to DMKCPI. WSCs =

0A - Mch Chk

11 - fatal I/O error

11 - Storage error (no attempt to recover from double-bit).

### Real Initialization Overview

All we have so far is a Nuc in storage. Init storage, work areas, control blocks, allocate save areas, mount + verify DASD, init directory, logon operator (now send msgs), start <sup>DMKDIS</sup> MIH, test for VMA + extended storage, allocate dump file, init error recovery, + wait for work.

DMKCPI + J control init.

↓

DMKSTA - init stor

PMA - init for a possible future PMA guest

API - AP/MP support if required (2 PAs)

MNT - Mount Devices - Search every RDEVBOOK  
Record status, volser, + if CP-owned.

SEG - init sys Vixt store (SPZINKs, etc.)

Pageable virtual store. Maps to seg + pag tables.

OPE - init Operator; looks for alt if req.

TOD - init clocks via prompt

Both proc. will synch on this

UDR - Build index into directory (pageable)

Build override file to chg class of cmds.

CPJ - finish init

## DMKCPJ



WRM - determine type of CP start

IDU - alloc dump

OPE - logon operator

CQN - spool File msg } simulate commands

CSO - start UR dev. }

SEG - init I/O to bring in the rest of pageable nucleus from sysres so it can be later paged out again to page space.

IOE - I/O err recovery

IOG - mch chk recovery

STA - Sys Pert Module: set disp timeslice + init monitor gathering.

OPE - fake STA MON cmd + AUTOLOG 1 as if the operator issued these commands.

DMKCPI - first thing it does is to set up the PSUs in PSA:

Restart - DMKPSADU

External - DMKEXTIN

SVC - DMKSVCIN

Program - DMKPRGIN - all lead to PRGnnnabend.

Machine - Disable wait 'B' at first, later will point to a valid entry point.

I/O - DMKIOTIN

Init PSA fields, CRs (0, 1, 2, 6, <sup>14, 15</sup>) clock comparator, CPU timer, CPUID, + TOD clock. Next issue STAP instr. to see if the proc. is an MP (store Address of Processor - Not installed unless machine is AP/MP capable). Init CPSTATUS fields in PSA



Find processor id of other processor (IPUADDRX).  
SIGP requires processor id; so each is tried.  
These run from 00-3F. It stops after  
finding one. All this goes into the  
absolute PSA. ARSTAT1 records this.

Determine VMA availability with E612 instr.  
that returns ucode level. This is not POO.  
Intercepts pgm chk while executing.

Shift RDEVBLOR to 3 bit shift if CP Assist  
is avail. + if less than 256K, CPSTAT5 tells  
3 or 4 bit shift.  
Call DMKSTA.

DMKSTA - find storage size. Use 55K  
to each 2048 byte area to set up keys  
up to 16 meg and watch for pgm check.  
Then, use TB instr to keep looking above  
16 meg. Store in DMKSYSRM.

Find if real size is larger than gen size to  
free unneeded core table entries.

Reserve two frames for PSAs. DMKPSA  
PREFIXA + B record ↑. Not yet filled in.

Test for 2K/4K storage keys. VM/56  
won't run w/ 4K keys. (CPSTAT4)

Allocate high free. Use DMKSYS if entered,  
or calculate default: Remember to omit V-R area.

pages/meg

< 16m 16pp/m - 1

> 16m 4pp/m

Multiply final result by 1.25 of hi free

{ Allocate prime storage or default (10%) alloc out of top end of high free.

Alloc trace table. < 16 4pp/meg (w/o V-R)  
> 16 1pp/meg

Set trace start, current, + end. (Also DMKPSAST)

Set Traceflg (+ '400') to 'FF' to indicate what events to trace. There are no commands to selectively trace. Just set the bits.

Discover what's left for DPA and init coretable. Freelist is built. Update coretable for CP pages as resident + CP owned. Call DMKCRET: will take end of DMKCPPE to page boundary to add to avail free storage.

Free list is built + ready (paging not yet happening).

DMKSTA returns to DMKCP1 after preparing storage for execution.

### DMKCPINT

Alloc areas in High Free. Preallocates CPEXBLOCKS. Point RUNUSER + LASTUSER to system + set up prefix registers.

### DMKAPI

SIGP xx STOP/START is used to control the other processor. A SIGP xx RESTART gets the other processor going after initializing the NIPU PSA loc 0 to an entry point in DMKAPI.

PSA  
PRIMEHI  
CO!  
In each PSA!  
They each have  
their own.

API copies absolute PSA to the other two locations and initializes the appropriate fields for each (mirror image) IPUADDR, LPUADDR, PREFIXA, IPUADDRX, LPUADDRX, and PREFIXB. Also, alloc save areas for 2nd processor. Init PXA, PXB for PSA extensions. Restart new PSW points to DMKAPIPR for AP. While API is running on the other side, the IPL processor loops (they are both using absolute PSA at this time). When API completes, start using prefixing and system lock for a double guarantee. Return to DMKCPI.

DMKAPIPR - Fixes restart new PSW back to dump routine. Check for CP assist on this processor (if not, neither can use it). Update APSTAT1. Find if system is AP or MP. If MP, change ARI OCT to point to appropriate chan table (since this processor has separate chan. set). A dev could be offline to one processor. Turn off spin of other processor + wait (DMKWAI).

DMKMNT - mount routine inits I/O subsystem:

get I/O lock since we have to run init routine on both processors. Touch all RDEVBLKs in order of chan table. Turn off disabled bit in RDEVBLK when dev answers. If never answered, then it stays disabled.

Flag alternate paths where discovered. (RDEVPTS)

Alloc RSPXBLOCK for UR devices.

(RDEVALT)

Find if Reserve / Release hardware exists.

If so, flag.

For DASD, read volser + detect duplicates.

Check if CP formatted (ownerid field in volume label = "CP370").

For CP owned volumes: find directory or override pointer in vol. label. If the volser is same as sysies, then store as primary. If not, these are alternates.

DMKALO - store ALOCBLOK pointer in RDEVBLOK.

Also, set CPOwned in RDEVBLOK.

DMKSEG - init CP segment, page, and subtables and core table. Check for VMA, EF (test protect...), + DAS (pgm call, pgm zfer) and record status. Check for extended key operation support and hardware segment protect. (Temporarily overlay pgm chk new psw).

DMKOPC - locate operator console. Call DMKTOD if not auto-starting. DMKTOD will synch multiple clocks if present.

DMKUDR - allocate CP vs for all entries in the director → userid, ACI group (RACF), DASD location of directory entry. Anchored at virtual address DMKSYSPL. (169 entries per yoro).

DMKCPJ - check to see if interval timer is running - set a value + loop while waiting for int. If no int, issue msg to turn on interval timer and loop again.

Determine type of start  $\rightarrow$  warm. Get type from operator. If response is shutdown, stop w/ WSC 6 + save warmstart data. Call DMKWARM.

DMKWARM -

Warmauto - read warmstart data, en termst lines. Put operator back on same console. Then join normal warmstart code.

Warm - read w/s data, + continue as if cold. DMRRSP is recreated as it was at shutdown. Also, find ACOCBLOCKs and point to the saved RECBLOCKs.

[Never warm start after changing the sysown list].

Init status of UR devices (classes, etc).

COLD - init chkpt cylinders. Chkpt spool + hold queues using DMKRSP anchors. If true cold, these will be empty. Re-init w/s cylinders.

CHKPT/Force - Restore spool files. Go through every spool file + recreate RECBLOCKs that can be hung off of ACOCBLOCK. If the SFBLK chains are broken, then a FORCE is required. Re-init w/s cylinders (fresh them).

DMKCPJ - set CPID to "CPCP" to indicate system is running. If it crashes warmstart may now be saved. Call DMKIDU.

DMKIDU - calculate needed dump space, build SFBLOK, RECBLK (DMKDMPSE), + allocate it. DMKDMPDV lists the device - dump, temp, or printer. Write symbol table (DMKSYM) to dump area - a primitive Nucmap. call VDG.

DMKVVG - build paging ACOCBLOKs.  
TDSK (clear TDSK w/DMKCPX).  
RECBLKs for PP + SW. (map shows all area avail. For spooling RECBLKs are not pre-built.  
Call DMKPSI to init paging storage. Also calls Vector Facility init.

DMKOPÉ - logon operator + issue messages

DMKCON - issue msg on sp files

DMKCSO - start UR

DMKHVD - init DMKKCPP bit map of VM/SP release level.

DMKSEG - bring in all pageable modules not in storage + page them out. - Go onto flushlist.

DMKIOE - init RMS. DMKMCHIN becomes mch chh new PSW.

DMKNLD - network load for 3705.

DMKSTPX - init scheduler info + start monitor.

RMS Init - Call DMKIOGFI. Issue STIDC to find avail channels. Determine if a chan exists w/o a RCHBLK. Build chan config table in DMKCH. Set CR15 to extended logout area. IOEL at X'AC' in PSA.

See if error recording cylinders are formatted (DMKSYS). They are formatted if not valid. Should always be valid.

Machines since 3081 do not do extended log outs & thus set a bit in the CPUID. CR15 will be  $\emptyset$ .

### Termination

DMKCP5 - handles shutdown. If AP switch to IPL processor & guresc. If MP,

set CP5HRCK in CPSTAT2 in both PSAs.

Stop prob state exec via CPINITD in APSTAT1 in both PSAs.

zero CR8 to stop MON. Close console file for operator & shutdown tube.

Record errors by chaining through RDEVBLKs (OBR records). Record DASD stats w/ DMKDSB.

Stop prefixing & goto dump routine

DMKOMP - enter at DMKDMPS & CRCP will be set in CPID field.

Abnormal - just about any pgm chk after initialization. The PIC is set in the PRG msg.

Abend macro - chain broken, DPA exhausted, etc.

SVC  $\emptyset$   $\rightarrow$  goto SVC DIE

DC CCB'ABC; ALI(N)

$\nwarrow$  abend code

PSW Restart: saves GPRs at X'500' +  
pass CNTR to DMKPSA X'0'. Aband  
code = PSA002. Goto DMKDMP.

enter  
at  
DK  
DMKDMP - detect if recursive. Save clock +  
regs to store status (software STSTS).  
Find if AP/MP, If yes, stop + store  
status on other processor. SIGP other  
proc. so it loads a wsc. Store status  
goes to absolute PSA, so move it to  
prefixed.

DMKDMPSEF is the  
SEBLOK  
Find dump device w/ DMKDMPDV. Default  
is printer.

CPID = 'CPCP'

If dump is to spool area, then  
things are a bit different. DMKSYM is already  
out there.

### DASD format of Dump

1 - DMKSYM

2-6 - reserved for now.

7-12 - storage pages 0-5

3-6 - storage keys built in storage pages 2-5

3-8-88

### Interrupt Handling

Gary

CP code always runs w/ I/O + EXT int disabled,  
in sup state, DAT off, EC on, + key 0.  
Subclasses of EXT Ints may be masked by  
bits in CR0.

### FCIHs

MCH - DMKMCHIN

I/O - DMKIOTIN

SVC - DMKSVCIN

Ext. - DMKEXTIN

Pgm - DMKPRGIN

Restart - DMKPSADU



Both CP & virtual machines can issue SVCs & have program checks. Therefore, to handle, CP checks the old PSW. If problem state then let the user handle it (except paging, of course). Since users run prob state, lots of privileged operation pgm checks occur for SCP guests. These get simulated (DMKPRV, DMKFPS) and not reflected to the guest.

The Dispatcher (DMKDSPWI) has an enabled window<sup>(s.m.s)</sup> right before dispatching a user to cut the probability of hitting an I/O'rapt as soon as user PSW is loaded.

The CP module DMKWAI loads an enabled wait & CP doesn't care if this is interrupted.

### Restart

Save regs @ DUMPSAVE

Set CPABEND to PSAPP2

Get dump dev

.. RECBLOK address

.. SFBLOK

Take Dump

Set CPID to 'CPCP' or 'WARM'

Simulate IPL

Machine Checks - 2 kinds: Recoverable, Unrecoverable.

Storage errors are usually recoverable. All others are usually not. Single bit errors are corrected by hardware & reported simply to log on EREP. Multibit errors

are not recoverable. CP can recover if it can get rid of the user + take the page frame offline (CORETABLE). If CP, the system abends. The user may be crashed even if the guest could handle it (like MVS). If the page w/ the error has not yet changed the page, then he is in luck; the bad page is abandoned, coretable marked bad, + the page is retrieved from aux store.

MCIC begins at X'ED.'

I/O Interrupt → DMKIOT (PMA interrupts go directly to guest).

3 possibilities: wait (Active wait), enable window, or user running. So, determine state. Save stat only if user running.

Check int dev addr +

located dev, cu, + chan blocks for dev. (DMKSCN)  
Switch timer to sys overhead

Check if dedicated chan + if so, pass directly to user. (DMKVIO)

RDEVBLK points to IOBLK. Requeue this to the dispatcher if I/O completed normally.

CE + PCI cause the IOBLK to be duplicated: one to dispatch + one to remain hanging on RDEVBLK.

- Unsolicited interrupts (no IOBLK off RDEVBLK)

- Deferred CC -

if ∅ do nothing, no need to 'pass up'

1 - DMKIOS - if Dev Busy wait + try again, if problem, reflect to user

2 - error processing - probably crash system

3 - DMKIOS

Unsolicited Int :

Create IOBLOCK + figure who handles

Some are ok - like hitting ENTER @ Terminal.

IOBIRA module

TTY - DMKCNS

Remote 3270 - DMKRG A

Sp Dev - DMKRSP

DASD - DASD

Local 3270 - DMKGRF

3704,05 - DMKRNH

Ded. Dev - DMKVZO

If an interrupt comes w/ unit check or from an unknown device CP will do a Sense.

(May simply not be in DMKRIO).

CP has no hot I/O recovery. It may build IOBLOCKs until mem is full or may thrash throwing the int away.

Errors :

If unit check, CP does sense before talking to user. When user finally does sense, CP pass the sense data. Record error ~~rebtly~~ to redrive z/0. Sense data goes into IOERBLOCK until passed to user.

↳ hangs off IOBLOCK + ends up off VDEVBLOCK.

CP does no recovery for virtual machine I/O. If it is to be done, the user must do it.

Chan check - DMKCCCH handles hardware related checks. Chan pgm chks, chan prot checks, + chan chaining chks are not hardware-related.

To free the controller.

For CCC, IFCC, CDC, try resetting the chan. + try again. A WSC 002 is loaded if the path is the last to a critical CP-owned dev + the reset fails.

DMKIOE handles CP I/O errors by examining the IOERBLK and the I/O extended logout.

Ext Interrupts - similar to I/O at the beginning.

Determine class of int (0 (350), 1-timer, 2-MSSF, AP/MP). Update count in PSA, cut trace entry, + put class + code in PSA ('x86' + '87')

	<u>Types</u>	<u>Action</u>
Class 0	0040 - Button	Disc System Oper.
	0080 - Real Interval Timer	DMKOSP
	0100 - Vint	Unit 201 Timer Assist ucode → DMKOSP
Class 1	1003 - TOD synch check	DMKCLK
	1004 - clock Comparator	DMKSTR
	1005 - CPU Timer - User's Whole Time limit Expired → DSP	
Class 2 AP/MP	2401 - MSSF (console) uses HCBLK	DMKMHC
	1200 - Malfunction Alert (MFA)	DMKMCT
	1201 - Emergency Signal SIGP	DMKEXTSL
	1202 - SIGP (External Call)	DMKEXTSL

XINTBLK are used to reflect an EXT Int to a virtual machine.

Try to keep running in UP mode. Cannot do if CP was running on the failed processor.

Single Processor Mode - CP has only one processor of a diatic. MVS will try SIGPs to other processor. Therefore, SIGPs are queued to the V=R user.

Other user has absolute PSA plus copies for each processors.

### SVC Interrupts

If <sup>prob state</sup> user issued SVC, then all are reflected to user except SVC 76 (logrec DMKVER) and SVC 179 (DMKCFM - adstop). Usually handled by VMA.

IF CP:

SVC 0 - abend

4 - reserved (abend)

24 - switch AP/MP - execution streams  
switch processors.

8 - Link-expansion of CALL macro.

12 - Return

16 - Release Save Area - Return to caller's caller

20 - Acquire Save Area - gets save area for me.

used for linkage among CP modules.

When adstop entered, an ADSTBLOCK is created + hung off VMBLOCK. Check address of SVC179 and look in ADSTBLOCK to see if it matched (+ code that was overlaid). If yes, stop execution. If no, reflect to user.

For SVC 76, verify that error needs recording + return control. MVS software recs appear invalid + therefore are passed back to the guest. Check GPRs 0 + 1 to see if DOS or OS, VS1 or CP. VM second level, uses SVC 76.

SVC saveareas: 12 DAs:

caller's return address R14

caller's base R12

R13

Save areas

Normally resident in prime storage.

Active savearea chaining: Savearea ↑ chain together.

Exit macro expands to:

LM R14,R11,SAVEAREA

SVC 12

Return address is placed in the PSW and the machine is permitted to run. Savearea is good for debugger. The running pgm, however can alter the saved registers.

SVC 24 - Put CPENBLOCK describing my current work on the other process queue (marked w/ affinity) + go look at the Dispatcher. The switch macro has a SIGP.

### Program Check Interrupts

MC causes a pgm check + is used regularly - the only normal CP pgm check.

DMKPRV to simulate.

DMKPRG - FCIH

DMKPTR - page from PIC 10 + 11 (TRESADD & 90' PSA)

DMKPAG - shadow tables

DMKPRV - PIC 2 (or 1 on occasion) DMKFPS is a front end.

Simulates any priv operation.

DMKVSI - virt I/O

VCM - " console

VSP - "

VMINST-  
instruction being  
worked on.

pageable

DMKIVC/D/E/F - Diagnose x'83xx

TMR - Timer operations (store, set)

DMKMPS - Virtual MP Mode - if in Single Processor Mode (SPM).

DMKEPS - duplicates other modules' functions but does it faster. Does 11 instructions.

Conditions: PER not active

instr on single resident page

virtual supervisor

virtual ecmode

PIC 13 - Special Op Exception  $\rightarrow$  for the semi-pri inst.

MVCP

MVCS

PC

SAC

SSAR

PT

result in a PIC reflected to the user (for guest simulation) to pretend the hardware doesn't support these.

PIC 80 - PEREvent  $\rightarrow$  DMKPERIL if for CP. If for user, save PER code + addr + ECBLCK, and notify user.

When reflecting pgm check, try to detect PILOOP. check pgm chk New psw + instr pointed to from it.

Privop processing must move the users pag down below 16meg for simulation.

DMKPRV has counters for each type of privop. VMINST holds instr while being worked on (+ may be residual).

3-9-88

Gary

## Real I/O Processing

PSA pointers:

Entry Length

ARIOCT → Real Chan Index Table

Holds displacement into Chan Blok Table for each chan (since they need not be coded in order). X'FFFF' means chan does not exist.

ARIOCH → Real Chan Blok Table

X'80'

Holds Table of all possible control units with displacement.

ARIOCU → Real CU Table

X'50'

Compressed displacement (4 bits) into Dev Table.

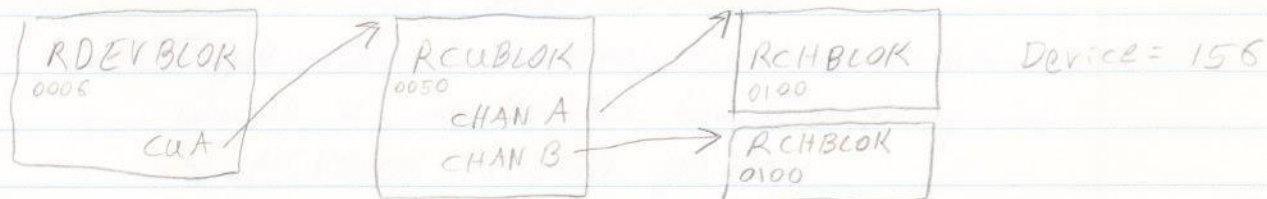
ARIODV → RDEVBLK Table

X'90'

For MPs, two sets of channels may exist. Since only 32 chans will fit on each processor, CP will build a second chan table. The other PSA will have a unique Chan Index Table.

The two Real Chan Blok Tables are contiguous. The same CU + RDEV blocks are used.

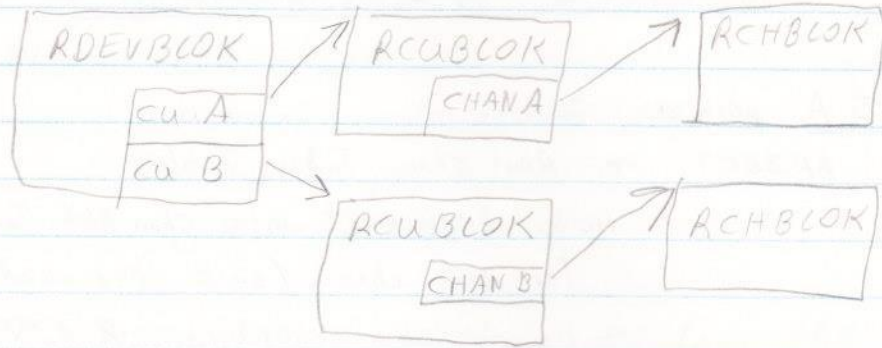
Each proc. has discovered during init. which CUs + Devs it can access.





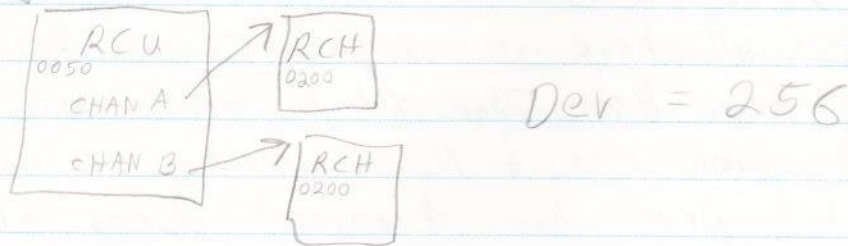
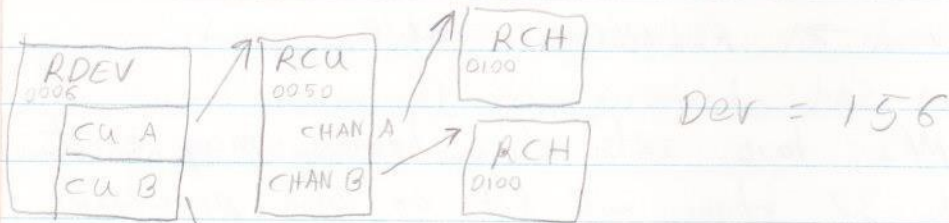
The two paths must lead to the same device.

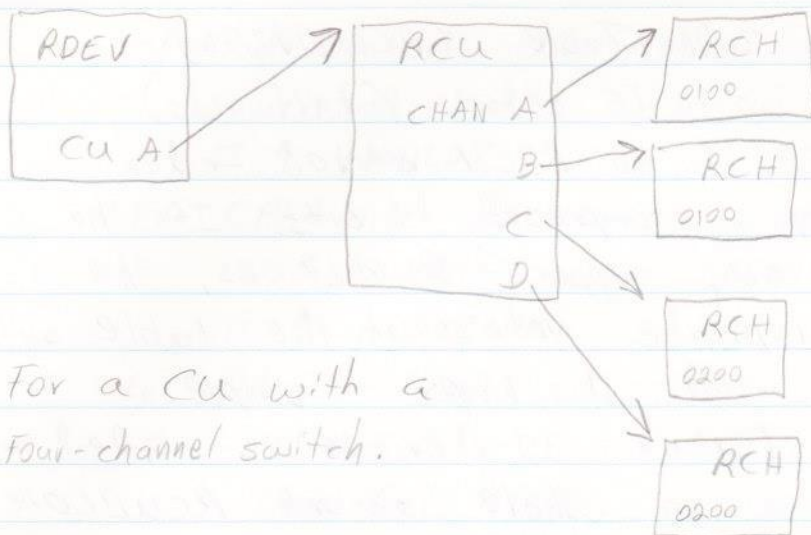
On an MP gen, an ALT CU may also be specified with the same address:



Note that the CHAN ↑ from the ACT CUBLOK is the 'B' field. This lets CP know that I/O may continue in event of a CU failure.

The ALT CU can be at a different address. This works on a UP too.





For a CU with a  
Four-channel switch.

All pointers are fullword pointers & are set  
to  $\emptyset - \emptyset$  when null.

RCHBLOK - the result of RCHAN macro.

Holds:

Address in form  $xx\emptyset\emptyset$

Owning Processor -  $0040, 0043$  for example.

Status - chan busy, offline, dedicated (if ded has VMS19)

Type - Selector, Byte, Blk mux

↑s to IOBLOKs queued on chan busy (RCHFTOB, LIOB)

CU index Table (32 entries, halfword each)

00 08 10 18

20 28 . . .

If dev is greater than 18 & the 18 entry  
is 'FFFF' back up on entry - to the 10 entry.

RCUBLOK - anchor in PSA ARIOCU

Address in form  $\emptyset\emptyset x \emptyset$  or  $\emptyset\emptyset x 8$

Status - Busy, Online, Scheduled

Type - Shared Sub-channels, Subordinate RCUBLOK

Queue of IOBLOKs for CU busy

4 RCHBLOK fullword pointers

## Device Index Table (RCUDVTBL) -

up to 16 entries (halfwords).

FFFF if no RDEV macro. Index entry is compressed 4 bits. If the CU only supports 8 devices, the appropriate entries in the table will be used. This keeps it simple.

If feature = 32 devices is coded on macro, a whole second RCUBLOCK is built with the address of the second string in the address field.

This 2nd RCUBLOCK has a RCUPRIME pointer & the subordinate RCUBLOCK bit is on in the TYPE field.

(instead of RCUCHA field)

The busy field in the first is the one that counts. The 2nd has 8 pointers to channels. All IOBLOCKs will be queued on the first.

An RDEVBLOCK will point to whichever RCUBLOCK is needed.

## RDEVBLOCK - anchored by ARIODV

Address → 000x

Type, Class (Data Areas, Appendix A)

Device Name - 3380, etc.

Status: RDEVSTAT + STA2, 3, + 4

Busy, scheduled, offline, dedicated

↳ can mean that SIO came back busy.

That cu will then owe a DE when free.

↑ First & Last IOBLOCK queued for device busy

↑ to active IOBLOCK (RDEVAIOB)

↑s to RCUBLOCKS (A + B)

Path Status (8 bits) 0 - avail

CU B  
xxxx  
CU A

For dedicated devices:

virtual addr

↑ VMBLOK

For CP-owned devices (in same CB location):

sysowned index (RDEVCODE)

↑ ALOCBLOK

overlay

For DASD:

volser

Last Cyl access

# links to this device (Q system)

Schedule Bit - when IOBLOKs arrive, they check first RDEV then RCU, then RCHAN.

At whatever level a busy is found, the IOBLOK is queued. When queued, it sets the schedule bit on the lower levels.

A second IOBLOK would spot the schedule bit + queue on that level. It knows that somebody is queued on need dev or CU + that a wait is required. When dev is free, I/Os are de-queued in Chan, CU, + Dev order.

IOBLOK (can't format w/ Analyze)

When active, ↑ from RDEVBLOK. Can queue behind Chan, CU, or Dev. When I/O done is queued to dispatcher w/ addr of routine to handle interrupt (SCIH). Holds: VMBLOK, virt dev addr, ↑ to real chan pgm (CAW), (IOBCYL) cylinder for first seek, + if queued: chain pointers (F+B).

After I/O interrupt, holds CSW + pointer to sense data if stored. (IOCRBLOK)

↑ than not for user

DMK2PRQ

IOBQDTC

Dispatcher IOBLOCK queue is usually empty since these are high priority tasks.

I/O comes from either CP or user.

CP - Paging

Spooling

Console Functions

Dir Maint

Error recording

SPTAPE

Monitor (if tape used)

Dump

CP must create the channel pgms.

Virtual Machine - prepares own chan pgms.

SIO to dedicated dev or to minidisk

Diag I/O

SIO to virtual spooling dev

SIO to " console

SIO to " CTCA

CP must be able to convert these channel pgms (+ for CTCA, simulate hardware).

Processing Flow

DMKIOS - the I/O Supervisor. Needs IOBLOCK as input (+ channel pgm).

DMKVSI - Virt SIO re-formatted + create IOBLOCK

DMKDGD - Diagnose I/O

DMKGRF - Graphics (for local)

DMKPAG - Paging

DMKRSP - Real Spooling

DMKIQ - Extension of DMKIOS to queue I/O

DMKIOT - Interrupt processor - updates IOBLOCK w/CSW + passes to dispatcher.

The dispatcher then passes control to the appropriate I/O SLIH - RSP, PAG, etc.

DMKIOS has two entry points:

QV - virt machine I/O (real addr not yet in IOBLOK)

QR - CP I/O (IOBFLAG indicates CP request)

R10 → IOBLOK

R11 → VMBLOK (stored in IOBWER)

R8 → RDEVBLOK or VDEVBLOK.

### Queuing Rules:

Q IOBLOK from lowest Busy component

Mark any lower components as scheduled

Sched component treated as busy for subsequent requests.

For ALT paths:

Search primary path first

On MP's, search current processor first

Queue on all busy paths (copy IOBLOK + chg dev adr.)

link IOBLOKs together - circular links

When avail path is found, dequeue other IOBLOK.

The duplicated IOBLOK are smaller → mini-IOBLOK. x'18'

Has only dev addr (for whatever path), flags, <sup>IOBFLAG bit</sup>

IOBLINK to original IOBLOK, chain points for other queued IOBLOKs, & ↑ to I/O component BLOK (IOBMQDIO).

Mini-IOBLOKs are pre-allocated at Init & limit the number of ALT path queue requests. Max is something like 15. DMKIOSMQ is ↑ to high free area where avail mini-BLOKs live.

When a real I/O is necessary, the addr from the Mini is merged into the real IOBLOK & the other mini BLOKs are discarded.

Find a path to a dev (DMKIOQ):

Check RDEVBLK

Check CU A then CU B

Check chans A-D

No paths; then return to caller's caller. (svc 15)  
DMKIOS's

Queuing -

- IOBLK queued off first unavailable path.
- Mini-BLKs queued off other paths.
- Search always starts at primary address.

Queuing order: on Dev:

- Get IOBLK Lock
- if not DASD, FIFO
- if fixed head paging - queue before moveable head requests
- if moveable head DASD - queue in ascending seek cyl adr order.

If CP, DIAG, or guest SIOF then issue an actual SIOF. If the guest issues a SIO then issue an actual SIO because guest would not be expecting a deferred CC.

If CC=3 to SIO(F) then mark path offline + try all paths (since only cu may be down). After all path tried unsuccessfully, mark Dev offline.

Moveable head DASD IOBLKs are dequeued in alternating ascending + descending order (as per RDEVSKUP).

DMRDIO - delinquent Int. Detection. DID sets a flag in RDEVBLK that is turned off by int processing.  $\rightarrow$  if dev found busy or sch. Thus, when DID comes back it can see if the device has stayed busy too long. If the dev is ok, then Busy will be on but DID bit will be off. DID then sets DID bit back on + will check again later.

If both Busy + DID bit on, then int. delinquent:

Attempt HDV (halt device)

Simulate IFCC (and software may be able to redrive)

If device is sch + no active I/O, then simulate cu end.

If busy bit on + no active I/O, then simulate device end.

Channel End is not simulated.

Always set MIT time longer for second level than for first level.

Gary

## Virtual I/O Processing

VMBLOCK anchors: VMCHTBL is index into VCHBLOCKS. 16 entries.

VMCHSTART  $\rightarrow$  Virtual Channel Blocks (VCHBLOCKS) <sup>2'30'</sup>

In numeric order w/ gaps for undefined channels. <sup>129'</sup> (2nd level is limited to 16 chans).

VMCUSTRT  $\rightarrow$  Virt. CUBLOCK for each 16 devices

No options for 8 or 32 devices.

VMDVSTRT  $\rightarrow$  VDEVBLK (2'30' length)



## VCHBLOK

chan addr  $\emptyset \times \emptyset \emptyset$

Status - busy, int. pending, dedicated

Type - selector, Block mux (Only  $\emptyset$  can be byte + is assumed to be so). The type of devices on chan determine type.

↑ RCHBLOK if chan is dedicated.

Virtual CU Index Table (16 entries)

CU Int. Pending bit map (halfword), 1 = pending.

## VCUBLOK - Anchored off VMBLOK

addr  $\emptyset \emptyset \times \emptyset$

Dev rupt bit map (halfword), 1 = pending

Status - busy or int pending

Type - Shared subchannels - determined by type of devices

Virtual Device Index Table

VDEVBLOK - all alloc in free storage in a contiguous chunk. Each time a dev is defined all are copied into a larger space. Created by ATT, DEFINE, + LINK. The request for free storage is conditional + if not available the request is denied. (ATT fails).

Dev addr =  $\emptyset \emptyset \emptyset \times$

Dev Type

↑ VMBLOK

I/O Count

↑ active IOBLOK (after SIOF) (VDEVIOB)

virt CSW (VDEVCSW)

Status - busy, rupt pending, not ready, dedicated

↑ RDEVBLOK if ded, dialed, minrdisk

Device Dependent Fields: console, spooled UR, + if minrdisk → DASD <sup>VDEVRCEN</sup> cyl relocation, DASD size (VDEVBN0),

in VMBLOK

VCHINT chan pending

and link chain pointer if LINKed to somebody's VDEVBLK. All linked VDEVBLKs are chained. The real owner's VDEVBLK may not even exist for the device.

### I/O Processing

The virtual machine puts a CAW into his page 0 pointing to the chan pgm in a virtual page:

SEEK

SET SECTOR

SEARCH ID Equal cccc HHHH RR  
0003 0011 07 ← a label not necessarily a count

TIC

READ - will put data into a virtual

page or pages.

Then the user issues a SIO or SIOF + does BC. Since these are Privops, CP finally knows what the virtual machine is trying to do.

Pgm Chk occurs. <sup>DMRPRG</sup> → DMRPRIV passes control to DMKVS1. It checks to see if virtual device is valid + if so, tries to turn the request into a real I/O.

DMKVS1 gets an IOBLK, translates virt dev addr to real, changes cyl to point within minidisk (Chan pgm is copied into RCW TASK where corrections occur). SET SECTOR is ok. SEARCH ID Equal is also left alone (except that real addresses are used for all of these). The READ addr is changed + IDAWs used. IDAW bit is turned on. IDAW entries are incremented on halfword boundaries.

The RCWTASK block is a fixed size + if the chan pgm is too long DMK VSI has to go get another larger RCWTASK block and copy everything into it + then continue until it fills again. Also, CCW translation continues until done or an error occurs + nobody else runs. Some products benefit by generating huge channel pgms (FDR + Syncsort)

The IOBLOK IRA points to DMKVIO.

CCWs are not translated for V=R user with SET NOTRANS ON using ded dev, chan or ATT devices.

DIAG 98 does not check at all.

For TIO or CURIO, the virtual BLOKs are checked. No I/O takes place (except CTCA). The CC is faked. HIO + HDV is handled similarly unless the virt blk's show a real I/O is in progress in which case CP will try to stop the real I/O.

### Completion:

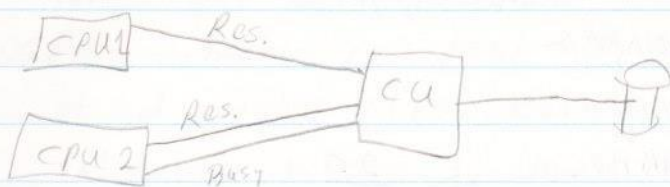
DMKVIO (SLIH) gets control from IOBLOK IRA, which checks CC (for possible deferred CC). If  $CC = 0$ , store CSW in VDEVBLK, + untranslate CSW address. If  $CC = 1$ , post pending Int. in VM I/O Bloks. If IOERBLK exists, store ↑ in VDEVBLK, set unit check in CSW.

$CC = 3$  not a path. Mark virt I/O blocks not busy.

DMKCCW is very large & complicated. (This is OCO in VM/SP XA).

### Reserve/Release Support

Hardware support of this is available on Tape and DASD. The lock is for a device & is held in the str. cnt. & blocks all but one path regardless of the # of CPUs.



CP does not use it, but it checks for hardware availability at init. For ALT paths, CP will not permit a user to issue Reserve. changes it to sense cc=0

Two MVSs under a single CPU are not protected from each other since both use the same path.

CP provides virtual Reserve/Release support.

Based on linking to a minidisk (which may be a full-pack minidisk). In directory entry MDISK statement MWV ALL ALL

↳ indicates virt. Res/Rel.

User may issue Res & CP will set indicator in VRRBLOK (anchored off VDEVBLOK).

↳ Holds

↑ VMBLOK of reserver

Virt. addr of owner of VDEVBLOK

CPEXBLOK anchor for everybody who got a busy & will eventually get a device end.

## Diagnose I/O

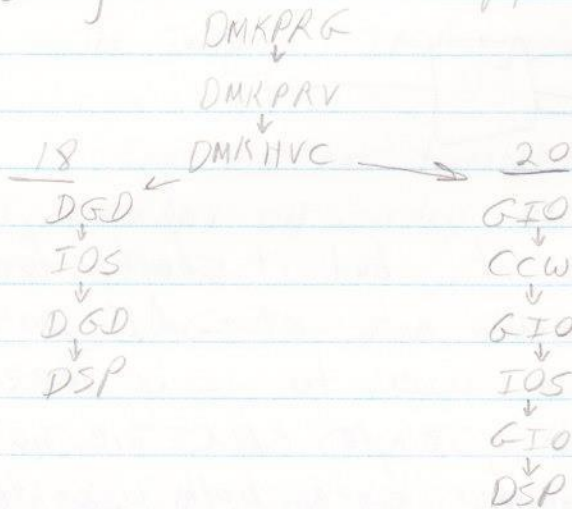
18 → Standard DASD I/O

20 → general I/O

Designed to run I/O synchronously. Wait until I/O is done + then deal w/ CC. Don't use virtual interrupt mechanism.

chan pgm must still be built.

Drag 83xx causes pgm chk.



DMKDGD - designed w/ CMS in mind. byte count must be 800/1024/2048/4096

CP will handle error recovery + user built chan pgm doesn't have to do it. User gets only a reg full of sense (rearranged).

Diag 98 - determine Real store addresses and avoid CCW translation. Subfunctions:

0 - give adv + lock pg

8 - DMKIOS executes, no translation

4 - unlock page frames

User must have DIA698 privilege specified in directory entry.

## Dispatcher

Dispatcher does more than dispatch. It also reflects interrupts to virtual machines.

3-10-88

Gary

DMKDSPCH - main entry point, requires global system lock.

If lock not avail, go to DMKDSPRU - the run user entry which cannot dispatch system work. Not used on UP.

DMKDSPA - fast entry for redispach after simulation or interrupt handling. If same user can't run, go back to main.

DMKDSPB - checks new PSW if changed outside of dispatcher (like CPSW in virt machine) for validation + exit to main CPA.

DMKDSPWI - enablement window prior to dispatch.

R11 points to VMBLOCK at all entries. This is to handle time accounting, reflect int., + validate virtual PSW. The scheduler may be called to modify the status of the virtual machine (it controls placement on the lists).

Dispatch order:

- IOBLOCKs (IOBIRA)
- TRQBLOCKs (TRQIRA)
- Enable for interrupts
- CPEXBLOCKs
- Tiny enablement window
- User
- Load enabled wait (or active wait)

If work is found at any level, it is processed + then the selection process begins again from the beginning.

## Reflect Interrupt

Restart - PMA  
Pgm  
Ext →  
Restart - all others

If user is dispatchable, select highest priority interrupt, verify virtual PSW enabled for interrupts, and exit if not enable. If ok, reflect interrupt by performing virtual PSW swap, validate new PSW, + store int code + related data in PSW.

Validation checks for PILoops + External Int loops.

If ccmode set on for user, build shadow tables in case an cc mode PSW is loaded. Once new PSW is determined, see if any more interrupts can be reflected (the new PSW may be enabled for more).

XINTBLOCKS are anchored in VMBLOCK at VMPXINT. Ext rrupts are chained in priority order.

I/O rrupts are indicated at VMCHINT with a bit map (lowest # chan interrupt is presented first). Depends on CR2 or BC PSW mask.

VMPEND indicates Restart Int pending.

VMRSTAT has all the bits that make a machine eligible to be dropped from the run list. This is usually idle wait for CMS.

The scheduler (DMK SCHED) alters dispatching status due to a change in the virt machine.

If Extends exist, only do work to fix → Process only paging I/O Blocks (IBLOCKS have real adr in first word. High order bit is off. TRQIBLOCKS

have the TOD is the first word & the high order bit is on). Highest dev addr is 1FFF. Process paging CPEXBLOCKS (identify by CPEXADD against table of known paging EPAs). NO users get dispatched during EXTENDS.

If no IOBLOCK, TRQBLOCKS, or interrupts, try CPEXBLOCKS. These are CP work that can wait a bit: page-in, error recording, logoff... So, load R15 from CPEXADD + R0-14 from CPEXBLOCK + then FRET the CPEXBLOCK. No trace entry except FRET of CPEXBLOCK.

Virtual Machine - check True Run List (TRL). Select first runnable user (TRL Lock + VMBLOCK lock). Recreate environment:

CR1 Seg Tab or Shadow Table

mode assists

PER if needed

FR Regs

GPRs

Create PSW (from old PSW or massaged VMPSW)

This code is complicated.

Insure CP stays in control of the hardware:

Load CRs with CP values.

Load interval timer with Dispatcher Time Slice.

Load CPU Timer with In-queue Time Slice.

CR0 - en all ext int. User can choose page size.

CR1 - seg Tab or shadow tables created by CP

CR2 - enable all chans

CR6 - set appropriate mode assists

CR8 - if CPTrop - CP loads



CR9-11 - CP PER or  
VM PER

CR14-15 - CPs Mch Chk setting.

Other CRs are not loaded. No DAS.

DMKWAI(ST) - steal VMBLOCKs from other processor. Get TRL lock + look for 2 or more dispatchable VMBLOCKs. Move it back to our TRL + release other CPUs TRL Lock. Then exit to dsp + discover something on TRL.

Both TRLs are usually equal in length & little stealing usually occurs.

Dispatcher Timeslice (DMKOSPQS) is the interval timer value. If compute bound (Q3) multiply by 4.

If no users in TRL who are runnable, see if users in runlist. If no, set CP idle wait. If yes, check VMRSTAT of each. If:

the sum of their ws is  $> \frac{1}{2}$  the DPA  
Then, set CP Pagewait!

Else, if one user is in I/O wait,  
then set CP IOWait

Else, go back to PAGEWAIT anyway!

Load CP Q3QD wait (on UP) or DMKWAITA for active wait. Switch CP PSCW to key 3 to indicate active wait.

For HPO running second level under VM/XA in MP mode, a DIAG 44 (XA diag!) is issued to tell XA that HPO would be running active wait and don't run until an interrupt comes in for HPO.

## Scheduler

Scheduler orders + maintains TRL which the dispatcher uses blindly. Objective: don't overload DPA. The TRL machines are entitled to have pages in storage. Runlist machines are eligible to have pages in storage. Favor interactive users over non-interactive users. Maintains data on individual users, groups of users, + DPA.

### Functions:

- Maintain lists
- Calc priority
- init Q-drop processing
- Determine if user is Q1
- Control E to runlist promotion
- Control Q3 users

### Lists:

Cyclic List - all VMBLOCKS: Sch doesn't care.

Eligible List - users waiting for storage. anchored at DMKSCHER.

Run List - DMKSCHRL - user who have storage allocated. User waiting on CP-simulation, etc.

True Run List - 1/CPU. Ready to run. Anchored in TRLANCHR in PXA + PXB.

TRLCT(+180) is count of VMBLOCKS on TRL.

Runlist } big slice  
Sch }  
Dispatcher - little  
↳ DMKOSPQS

VMBLOCKs are ordered by deadline on queues. Access is not necessarily round-robin. Queue Manager Units (QMUs) are bits 5-38 of TOD clock and is the deadline. Take TOD now + add average cycle time through eligible + run lists and make that the deadline (subject to coefficients).

### Time Slicing

Quantum - determined by processor speed (determined by BC loop) in DMKSTP. This is stored as a word @ DMKOSPQS + is the length of a single stay on the processor. This value is loaded into the interval timer.

The dispatcher time slice is how long a machine can stay in queue:

$$Q1 - 8 * \text{Quantum} \quad (\text{DMKSCH} + Q1) \downarrow \text{VMQBLOK}$$

$$Q2 - 64 * \text{Quantum} \quad (\text{DMKSCH} + Q2) \uparrow \text{anchors}$$

Quantum can be changed by SET SRM DSPSLICE

If a user has not entered wait by end of dispatcher ts, mark user compute bound + get a worse position in TRL. If CPU is not scarce, maybe increase timeslice value to cut number of trips through the scheduler.

If Runlist time expires w/o relinquishing control, you get Q2 (non-inter) + are dropped from TRL + run list. Logical susp out is initiated. Priority is recalculated (VMEPRIO) to do next deadline. Storage requirements are also recalculated → a guess at best. (VMPROJWS) Since a user is coming off TRL, see if another can go from E list to TRL.

When user goes idle, drop from TRL + set 300ms timer. At pop, drop from runlist + promote somebody else. For dropped user, recalc deadline priority + storage requirement. If user comes ready before 300ms put back in TRL + let him use runlist time slice (+ get downgraded to Q2).

A user stays Q1 by entering a wait longer than 300ms before the Runlist timeslice expires.

The deadline TOD starts at zero for moment of IPL. This gives 2 yrs of bits.

$VMEMPRIOR$  - E list priority

Equals  $VMTPRIOR$  except for compute bound machines.

$$VMEMPRIOR = VMQPRIOR + \underbrace{TODCURR - TODIPL}_{\text{this "normalizes" TOD}}$$

$VMQPRIOR$  is the delay factor.

Directory priority value undergoes conversion:

$$\text{Factor} = \frac{\text{Factor}/64}{64}$$

A lower priority results in better service. See the CP logic Manual for a complete description.

CPU + Paging load factors also count. The individual user perf. is \* by average perf. If user is better than average then product is fractional + priority becomes a lower number.

If a user comes ready after a long wait, a new deadline is calculated which improves position in E list.

Projected PWS: depends on # pages + CPU time. How much work was done per page read. It favors locality of reference. This is compared to ideal. PWS is adjusted w/in bounds,  $PWS < \text{max pages used}$ , and  $PWS > \frac{1}{2} \text{ average \# resident pages}$ . Shared pages don't count in PWS.

When moving from E list to runlist:

- Check 1st Q1 user on E list
- If PWS fits, move to runlist
- Are enough Q1 users in runlist? If not add anyway.
- Can a Q3 on runlist be pre-empted?
- See in non-inter user can fit.

Criteria are kept in VMSBLK.

Min/Max % avail DPA

Min # users in runlist = 1

# users in runlist

Min/Max avail DPA pages

Sum of PWSs of runlist users.

Priority for non-interactive:

When user used 6 consecutive Q2 timeslices - Q3 flag is set and will get a terrible deadline but will get 8 Q2 timeslices when it finally runs. The stay on the E list will be 8 times longer.

## Commands

### Priority Calc

SET Priority  
SET Favored %  
SET SRM PB - set high for paging problems.  
SET SRM IB - unnoticable  
SET SRM DSPSLICE - maybe larger if no CPU bottleneck.  
SET QDROP OFF - won't be swapped.  
SET Paging % -

### E'list to Runlist Promotion (PWS calc)

V=R option - PWS =  $\emptyset$   
LOCK } subtracted off WSS  
SET Reserve }  
SET SRM MAXWSS - can cause thrashing  
SET SRM ZBUFF - should be lower than the default 50%  
Favored

### To Monitor the Scheduler:

Ind Users

Ind Queues } Reports detailed scheduling info.  $\rightarrow$  VMEPRIOR for users.  
Ind Position }

Q - anything SET

DMKSCHQ1 - timeslice (Dw) - 1st entry in VMQBLOK

DMKSCHQ2 -

Each has # users in runlist \*

Only one E'list

# users in E'list (Anchor of E'list is DMKSCHED)

The VMSBLOK has slightly different categorizations:

- Interactive
- Batch

and the counts will not agree

Gary

## VM Assists

These are extensions to the 370 P00 and are dependent on processor models. These do not have to be available.

VM can still run (except PMA). Assists include:

Virtual Machine Assist

ECPS: VM/370 - includes CP Assist

MVS Extensions - supported for MVS second-level

PMA

Virt Machine Assists: reduces privop simulation. Documented in 1980 manual: GA22-7074

ECPS - smaller processors only (4300s). Some CP code ran as microcode. Does not aid the user (i.e. Privop), helps CP - find freestorage, etc.

Expanded VMA - added more Privops to ucode.

Virt Timer Assist - replaces CP code on updating Virtual Interval Timer. - Not implemented on larger CPUs

MVS Extensions - DAS, LAP, & common segment facility.<sup>PC</sup> HPO implements with additional ucode, simulation, & reflection of non-availability. Documented in HPO 3.6 System Programmer's Guide.

PMA - avail only to V=R user. Designed to help MVS guest. Permits true supervisor state for MVS, & I/O rpts handled by MVS through new interfaces between CP & guest. CP overhead is reduced.

### VMA

Availability is determined by DMKCP1 at init time. Used for all users → SVC may be reflected directly to guests. Also, improves shadow table maintenance. Default is to be used by all machines but may be turned off. CR6 controls VMA. Points to microcode block.

## Flow:

User in virt supervisor + issues privop.  
Microcode detects before pgm chk + checks CR6  
to see if supported. If simulation possible,  
do it + resume user execution. User is  
not interrupted, traced, or monitored. CP  
doesn't even know a privop was issued.  
User prob state time increases.

## CR6:

bits 0-7 → 0 VMA enabled

4 → 1 SVC assist disabled

2 → Extend key operations

3 → ECMode Setting

8-28 → ↑ MICBLOK (in high free) <sup>1/machine using any assists</sup>

↑ SEGTABLE

↑ ECBLOK

Flag | ↑ VMPSW

↑ work area

↑ virt interval timer

Flags

29 - 370 EF enabled

30 - PMA enabled

1 - Prob/Super State

## Instructions

IPK, SPKA

ISK, SSK, RRB

SSM

LRA

LPSW

STCTL

STNSM, STOSM

SPT, STPT

ISKE, SSKE, RRBE

All SVCs except 76  
are directly reflected to  
the user.

Shadow Table - prevents PIC 11  
if the shadow table only  
need be validated with  
data from CP tables.



Assists will be disabled if user is tracing, using PER, or has an ADSTOP set. HPO running second level cannot offer assists to its users.

### CP Assist

Certain CP functions are now in ucode. Initiated via EB instruction:

Get/Ret Free Stor

Loc I/O CBs

lock/unlock page

DSP returns

translate CCWs

SV8 + 12 (Link + Ret)

Assists can be individually disabled. If successful, a trace entry is created with  $\chi'80'$  ORed on in first byte to alert that CP Assist did it. If function fails, CP performs through normal code.

At init, if not supported, the EB instructions are all overlaid with no-ops. The assist may not exactly match the original CP function.

CP Assist is not implemented on the big CPUs because there is no speed gain.

### Expanded VMA

Gets control is regular VMA can't do a function.

LPSW

SCKC

STNSM, STOSM, SSM

TCH

SIO (SIOF)

DIAG

SPT, STPT

