

# XvMotion overview

# Migrate\_ToStart

MigrateVigorStart

    Migrate\_ToStart

        MigrateRPC\_MigrateBegin

        MigratePlatformInitMigration

            vmxParams.hasDiskPreCopy = SVMotion\_IsSVMotion()

            VMKernel\_InitMigration(&tmpSpec, &vmxParams)

            MigrateTypeBridge\_InitMigration

                migrateTypeBridge[mi->type].fns.\_name(...);

// MIGRATE\_TYPE\_VMOTION for vmotion and TYPE\_FSR for fsr

# VMotion\_InitMigration

```
VMotion_InitMigration // called for vmotion, xvmotion
    vi = Migrate_Alloc
    VMotionInitStream
    MigrateNet_SetNetStackInstance
    VMotionCreateHelperWorlds
    VMotionStartHelperWorlds
    if (VMotion_IsSource(vi))
        VMotionQueue_Add(vi, VMotionSend_ResolveFeatures)
        ...
        VMotionQueue_Add(vi, VMotion_ReportUserRPCState)
        VMotion_PostVmxMsg(vi, MIGRATE_VMKMSG_ENABLE_USER_RPCS)
    VMotionQueue_Add(vi, VMotionSend_ResolveSwapType)
    VMotionQueue_Add(vi, VMotionSend_ResolveRtt)
```

# FSR\_InitMigration

```
FSR_InitMigration // called for svmotion
    FSRAttachToRemoteFSR
        // Connects the two FSRInfos together. Checks that the memory
        // sizes and swap types match.
        FSRPostVmxMsg(fsr, MIGRATE_VMKMSG_ENABLE_USER_RPCS)
        FSRPostVmxMsg(remoteFSR, MIGRATE_VMKMSG_ENABLE_USER_RPCS)
    // start swap copy
```

# Migrate\_RPCsReady

really the start of everything

```
Migrate_RPCsReady // MIGRATE_EVENT_RPC_READY event
  MigrateInitiateCopy(&msgList) // no multiwriter
    MigratePlatformToStart
      // memory copy start for vmotion
      // setup storage streams for xvmotion
      // no-op for svmotion
    SVMotion_Start
      // copy for xvmotion and svmotion
      // no-op for vmotion
```

# MigratePlatformToStart

MigratePlatformToStart

```
hasMemPreCopy = !(migrationState.spec.type == MIGRATE_TYPE_FSR);  
hasDiskPreCopy = SVMotion_IsSVMotion();  
VMKernel_MigratePreCopyStart(hasMemPreCopy, hasDiskPreCopy)
```

FSR_PreCopyStart FSRTransferDVSSState FSRTransferDVSOOBRuntimeState
---

# MigratePlatformToStart

MigratePlatformToStart

```
hasMemPreCopy = !(migrationState.spec.type == MIGRATE_TYPE_FSR);  
hasDiskPreCopy = SVMotion_IsSVMotion();  
VMKernel_MigratePreCopyStart(hasMemPreCopy, hasDiskPreCopy)
```

VMotion\_PreCopyStart

if (!hasDiskPreCopy)

VMotionQueue\_Add(vi, VMotionSend\_PreCopyStart)

VMotionSend\_RemoteCall(vi, VMOTION\_MSG\_PRECOPY\_START)

VMotionStream\_StartAsync(VMOTION\_STREAM\_PRECOPY)

Action\_Post(vmmLeader, MONACTION\_MIG\_PRECOPY\_START); // monitor

else

XVMotion\_SetupMigration

xvm = (XVMotion \*)VMotion\_Alloc()

VMotionQueue\_Add(vi, XVMotion\_DiskPreCopyStart)

VMotionStream\_StartAsync(vi, VMOTION\_STREAM\_STORAGE)

# Migrate\_RPCsReady

really the start of everything

```
Migrate_RPCsReady // MIGRATE_EVENT_RPC_READY event
  MigrateInitiateCopy(&msgList) // no multiwriter
    MigratePlatformToStart
      // memory copy start for vmotion
      // setup storage streams for xvmotion
      // no-op for svmotion
    SVMotion_Start
      // copy for xvmotion and svmotion
      // no-op for vmotion
```





# SVMotionMirroredModeThread

```
SVMotionCreateDestDisks()  
    SVMotionDiskCreate  
        if (disk->isRemote)  
            SVMotionRemoteDiskCreate  
                SVMotionSyncRPC(MIGRPC_OP_SVMotionDiskCreate)  
                // on dest  
                SVMotion_DiskCreateRPC  
                DiskLib_Create  
        else  
            SVMotionLocalDiskCreate  
                DiskLib_Create  
SVMotionLoadDestDisks()  
SVMotionStun() // install mirrors  
SVMotionThreadWaitForStun()  
Vmkevent_RegisterHandler(VMKEVENT_GET_DISKCOPY_BITMAP,  
                        SVMotion_SendDiskCopyBitmap)  
SVMotionMirroredModeThreadDiskCopy
```

# SVMotionMirroredModeThreadDiskCopy

```
SVMotionMirroredModeThreadDiskCopy
```

```
    SVMotionInitDiskCopySchedulerQueue
```

```
        schedulerQueue[(*iter)++].disk = disk;
```

```
    SVMotionAsyncDiskCopySort
```

```
        schedulerQueue[i].weight to: 0=in-progress, -1=completed, n=to-go
```

```
    // 4 at a time, and if weight >= 0
```

```
        SVMotionAsyncCopyInit
```

```
            SVMotionMirroredModeAsyncCopyIoctl
```

```
                Mirror_SendAsyncCopyIoctl // SVM_ASYNC_DISK_COPY_START
```

# SVMAsyncDiskCopyStart

SVMAsyncDiskCopyStart

kernel land

// starts a SVMAsyncDiskCopyThread

SVMAsyncDiskCopyThread

for every 64MB

copyArgs->dm.fsdm.dmExtent[0].length = 64MB

copyArgs->dm.fsdm.dmExtent[0].srcFileOffset = progress->currOffset;

copyArgs->dm.fsdm.dmExtent[0].destFileOffset = progress->currOffset;

SVMFDSIoctlMoveData // not really an ioctl anymore

SVMDMCopyDiskUsingBitmap // XXX leave for later

SVMMoveData

SVMMoveData

if local and DM

FSS\_IoctlByFH(IOCTLCMD\_VMFS\_VMKMOVEDATA) // DM

if xvmotion

SVMReadDataToXVMBuffer

# SVMReadDataToXVMBuffer

API for xvmotion transfer of 64MB

SVMReadDataToXVMBuffer

    SVMXVMProcessVMFSBlockData // vmfs optimizations

for all dm->numExtents // usually 1

    SVMAsyncIORead // beware, xvmotion specific

        MigrateBridge\_XVMAllocBlocks

        SVMAsyncIOIssueRead

            FSS\_AsyncFileIO(comp=SVMAsyncIOReadDone)

                // on completion

                SVMAsyncIOReadDone

                MigrateBridge\_XVMEnqueue

    SVMWaitForReadIOCompletion

# Summary of [s/x]vmotion disk copy

- vmx sends an start copy ioctl to the kernel
- kernel starts a thread that copies 64MB chunks at a time
- for each chunk, call the DM or the xvmotion DM
- mirroring? next

# Mirroring setup

```
init_module
    SVMRegister
        FDS_RegisterDriverWithAttributes(SVM_DRIVER_NAME, &svmOps, driverID)
        svmDriverID = driverID;
```

```
static FDS_DeviceOps svmOps = {
    SVM_AsyncIO,
    SVM_ioctl,
    SVM_MakeDev,
    SVM_RemoveDev,
};
```

```
SVMotionDiskCreateMirrorNode
    VMKernel_FDSMakeDev
        SVM_MakeDev
            DevLib_CreateDevice(devLibObj, svmDriverID)
```

# SVM\_AsyncIO

```
SVM_AsyncIO(FDS_HandleID)
```

```
    if (ioFlags & FS_READ_OP)
```

```
        return FSS_AsyncFileIO(driverData->fhids[0]) // to the source only
```

```
    if (driverData->mode == SVMMIRROR_MODE_LOCAL)
```

```
        SVMAsyncIOLocalInt
```

```
    else
```

```
        SVMAsyncIORemoteInt
```



# SVMAsyncIOLocalInt

SVMAsyncIOLocalInt

```
if ((endOffset <= driverData->dmCopyStartOffset) // we already copied it
```

```
    Async_StartSplitIO // complete when everybody completes  
    for (index = 0; index < ARRAYSIZE(driverData->fhids); index++) {  
        FSS_AsyncFileIO(driverData->fhids[index])
```

```
else // source only
```

```
    FSS_AsyncFileIO(driverData->fhids[0]) // 0 is source
```

# SVMAsyncIORemoteInt

SVMAsyncIORemoteInt

```
if ((endOffset <= driverData->dmCopyStartOffset) // we already copied it
```

```
    MigrateBridge_XVMAllocBlocks
```

```
    MigrateBridge_XVMEnqueue // we already have the data to write
```

```
    FSS_AsyncFileIO(driverData->fhids[0]) // 0 is source
```

```
else // source only
```

```
    FSS_AsyncFileIO(driverData->fhids[0]) // 0 is source
```

# XVMotion buffer

MigrateBridge\_XVMAllocBlocks

    XVMotion\_AllocBlocks

        XVMotionCheckForBlockMemLocked // wait

        XVMotionAllocSgaBlocksLocked // alloc

MigrateBridge\_XVMEnqueue

    XVMotion\_EnqueueBlocks

        for (i = 0; i < sga->length; i++)

            XVMotionQueuePush(xvm->slice, dq, blk, offset, len)

# SVMotionMirroredModeThreadDiskCopy

```
SVMotionMirroredModeThreadDiskCopy
    // 4 at a time, and if weight>=0
        SVMotionAsyncCopyInit
            SVMotionMirroredModeAsyncCopyIoctl
                Mirror_SendAsyncCopyIoctl // SVM_ASYNC_DISK_COPY_START
            // wait for all copies to finish

SVMotionVMKernelPreCopyIterDone
    VMKernel_MigrateDiskPreCopyIterDone
        VMotion_DiskPreCopyIterDone
    VMotionQueue_Add(vi, VMotionSend_PreCopyStart)
        // start memory precopy
```