

**UNIVERSIDADE FEDERAL DE SANTA CATARINA
PROGRAMA DE PÓS-GRADUAÇÃO EM CIÊNCIA DA
COMPUTAÇÃO**

Marcelo T. Pereira

**RIFFS: Um Sistema de Arquivos para Memórias Flash
baseado em Árvores Reversas**

Trabalho individual submetido à Universidade Federal de Santa Catarina como parte dos requisitos para a obtenção do grau de Mestre em Ciência da Computação.

Orientador:
Antônio Augusto Medeiros Fröhlich

Florianópolis, Fev de 2004

RIFFS: Um Sistema de Arquivos para Memórias Flash baseado em Árvores Reversas

Marcelo T. Pereira

Esta Dissertação foi julgada adequada para a obtenção do título de Mestre em Ciência da Computação, área de concentração Sistemas Operacionais e aprovada em sua forma final pelo Programa de Pós-Graduação em Ciência da Computação.

Fernando Gauthier

Banca Examinadora

Antônio Augusto Medeiros Fröhlich

Marcelo Pasin

Rômulo Silva de Oliveira

Wolfgang Schröder-Preikschat

*“A melhor forma de prever o futuro é criá-lo.”
(Peter Drucker)*

“às minhas alianças afetivas,
à natureza,
ao futuro...”

Sumário

Lista de Figuras	viii
Lista de Tabelas	ix
Resumo	x
Abstract	xi
1 Introdução	1
2 Memórias Flash	4
2.1 Conceitos Gerais	4
2.2 Operações	5
2.3 Tecnologias	7
2.4 Estudos de Casos	8
3 Sistema de Arquivos	11
3.1 Dispositivo de Armazenamento	11
3.1.1 Blocos Lógicos	12
3.1.2 Gerenciamento de Blocos Livres	12
3.2 Gerenciamento de Arquivos	13
3.2.1 Operações	14
3.2.2 Gerenciamento dos Blocos de Arquivos	15
3.3 Gerenciamento de Diretórios	17
3.3.1 Operações	18

4	Sistemas Arquivos para Memórias Flash	20
4.1	Conceitos Gerais	21
4.1.1	Apagamento-e-escrita	21
4.1.2	Remapeamento	21
4.2	Device Drivers	24
4.2.1	Estudo de Casos	25
4.3	Sistemas de Arquivos	27
4.3.1	Estudo de Casos	28
5	Projeto do Sistema RIFFS	33
5.1	Motivação	33
5.2	Objetivos	35
5.3	Modelo Arquitetural	37
5.3.1	Sistema de Arquivos	37
5.3.2	Dispositivo Armazenamento	39
5.3.3	Gerenciamento de Diretórios:	42
6	Implementação do Sistema RIFFS	44
6.1	Componentes do Sistema	45
6.1.1	Flash Control	45
6.1.2	Scanner	46
6.1.3	Allocator	47
6.1.4	Device Manager	48
6.1.5	File Manager	48
6.1.6	Directory manager	50
6.1.7	Garbage Collector	51
6.1.8	Notas	53
6.2	Resultados do Sistema	53
6.2.1	Plataforma de testes	53
6.2.2	Escrita de Dados	54

6.2.3	Codificação do sistema	57
6.2.4	Tamanho de estruturas	57
7	Conclusão	60
	Referências Bibliográficas	63

Lista de Figuras

4.1	Atualização de dados na Flash.	22
4.2	Atualização de dados na Flash.	23
4.3	Estrutura da FTL.	26
4.4	Arquitetura do MTD.	28
4.5	Arquitetura do JFFS2.	31
4.6	Camada TrueFFS dentro do Sistema Operacional.	32
5.1	(a) Visão Lógica. (b) Visão da RAM. (c) Visão da Flash	38
5.2	Arquitetura do Setor.	40
5.3	(a) Visão Lógica. (b) Visão na RAM. (c) Visão na Flash	43
6.1	Módulos da arquitetura RIFFS.	44

Lista de Tabelas

6.1	Tempo de Escrita em Arquivo (TEA) para o RIFFS	56
6.2	Tempo de Escrita em Arquivo (TEA) para o JFFS2.	56
6.3	Comparação de desempenho entre: RIFFS e JFFS2.	57

Resumo

Este trabalho apresenta uma nova estrutura de armazenamento de dados em memórias flash, chamada de “Reverse-Indirect Flash File System” (RIFFS). As Flashs possuem uma limitação na atualização de seus dados, e pensando em amenizar esta característica pensou-se em deixar todos os dados e meta-dados dentro do próprio arquivo. Isso seria impraticável com os sistemas existentes, porque não seria possível localizar um arquivo diretamente, a partir do “nodo raiz” da árvore. A maneira encontrada foi criar uma árvore reversa. Este esquema quebraria a navegabilidade do sistema, e então uma árvore direta precisa ser construída na memória RAM. É mostrado neste trabalho o gerenciamento de uma árvore reversa para contornar as limitações da memória flash. Dentro deste esquema é possível evitar excessivas atualizações e operações de escrita, aumentando assim a vida útil da flash.

Keywords: sistemas operacionais, sistemas embutidos, memória flash, sistema de arquivos.

Abstract

This project presents a new technique for flash storage management called a “Reverse-Indirect Flash File System” (RIFFS). However, flash memories have a drawback: its data cannot be updated-in-place. To solve this limitation, all the data and meta-data is leaving inside of the proper archive. This would be impracticable with the actually systems, because it would not be possible to locate a file directly, from “root”. The solution was to construct a reverse-tree. This schema would break the navigability of the system, and then a direct tree need to be constructed in RAM memory. This work shows the reverse-tree management schema to solve the limitations of flash memories. This solution helped to minimizate extreme updates and write operations, increasing flash life-time.

Keywords: operating systems, embedded systems, flash memory, file systems.

Capítulo 1

Introdução

Já não é de hoje que o Homem moderno se acostumou com aparelhos ou ferramentas eletrônicas no cotidiano. Alguns trazem vantagens ao desempenharmos nosso trabalho, enquanto outros nos trazem conforto, comodidade, etc. Estes aparelhos estão presentes nas mais diversas áreas, desde a mais simples como um relógio despertador, até as mais complexas como por exemplo o sistema de foco de uma filmadora digital ou o sistema de navegação de um carro [GRO 96]. Para o funcionamento destes dispositivos, existe a necessidade de um sistema de configuração e controle, e a todo este conjunto damos o nome de **sistemas embutidos** ou **sistemas dedicados**.

Os sistemas embutidos são construídos com componentes eletrônicos em geral, como por exemplo circuitos integrados, portas lógicas, circuitos impressos, microprocessadores, etc; e comumente controlados por softwares específicos. Normalmente quando nos referimos à estes circuitos, lembramos de computadores pessoais (PC) e de seus “*chips*” de processamento (CPU), esquecendo dos vários equipamentos a nossa volta que também se utilizam deles. De acordo com Tennenhouse [TEN 00] apenas 2% dos 8 bilhões de processadores fabricados em 2000 foram aproveitados para estações de trabalho (PC), sendo que a grande maioria teve seu fim em sistemas embutidos.

O sucesso dos sistemas embutidos não se deu apenas à utilização da tecnologia de componentes eletrônicos, mas também ao uso de softwares específicos para cada sistema. Estes programas podem tanto exercer apenas algumas funções dentro do

sistema como executar funções mais complexas de gerenciamento de processos, alocação de recursos, etc.

Independente da complexidade e do tamanho, estes softwares específicos precisam ser armazenados em algum tipo de mídia não-volátil, e a mais comumente usada é a memória apenas de leitura (Read-Only Memory - ROM), e suas variantes: EPROM¹ e **EEPROM**². Este último tipo de memória foi eleita pela sua baixa latência de leitura, alta resistência e pequeno tamanho, em comparação à outras mídias. Quando é necessário atualizar algum dado, este tipo de memória precisa ser apagado por inteiro e depois reescrito. A fim de melhorar a rotina de atualização, foi agregado mais um membro à esta família de memórias, chamado de **Memória Flash**. Esta não precisa ser apagada por inteiro, mas em blocos chamados de **unidades de apagamento** ou **setores**. Com isso sua atualização acaba sendo mais rápido e conseqüentemente seu consumo de energia acaba sendo menor.

Com uma alta densidade, um peso leve, um pequeno tempo de latência, um baixo consumo de potência, e por fim uma vantagem na atualização de dados, as memórias flash se tornaram um meio atrativo de armazenamento de dados. Infelizmente por causa do seu preço, hoje em dia ela não é usada como armazenamento principal em computadores pessoais (ou outros sistemas de grande porte), mas suas vantagens de armazenamento em sistemas embutidos é clara.

Como dito anteriormente, existe o grupo de SE que está ficando cada vez mais complexo, como é o caso do telefone celular, por exemplo. Esta linha de dispositivos precisa se preocupar, entre outras coisas, com o armazenamento: dos dados de configuração, dos dados do usuário, de módulos do próprio sistema (como por exemplo uma nova versão da máquina virtual java), etc. Por causa do tamanho das aplicações, houve uma necessidade do aumento das memórias flashes, e da implementação de um sistema de arquivos. A manipulação de dados dentro destas memórias possui suas peculiaridades, e por este motivo, a construção de um sistema de arquivos em uma flash torna-se um desafio à engenheiros de software.

¹EPROM - Erasable ROM

²EEPROM - Electrically EPROM

Pensando em aproveitar as vantagens da flash, e tentando contornar suas limitações, pesquisadores tiveram que desenvolver e “reciclar” conceitos para construir sistemas de arquivos para essas memórias, manipulando os dados de uma forma eficiente.

Este trabalho apresenta uma estrutura de armazenamento diferente dos sistemas de arquivos para memórias flash atuais. O sistema de diretórios é baseado em **árvores reversas** e o sistema de arquivos é baseado em uma estrutura diferente das encontradas atualmente, chamado de **contexto de arquivo**.

A principal idéia deste artigo está nas estruturas físicas gravadas na Flash, e como gerenciá-las eficientemente. Neste cenário, cada arquivo agrupa todas informações pertencentes a ele próprio, criando assim o **contexto de arquivo**. Com isso é possível uma diminuição na atualização dos dados de controle em relação a outros sistemas de arquivos. Para conseguir satisfazer este requisito, foi preciso a concepção de uma **árvore reversa** com a ligação **indireta** entre seus nodos, surgindo assim o nome do projeto: Reverse-Indirect Flash File System - RIFFS.

O próximo capítulo mostra as desvantagens das memórias flash, e possíveis algoritmos para contorná-los. O terceiro capítulo é um estudo sobre as teorias clássicas de sistemas de arquivos, e seus algoritmos. O quarto capítulo mostra como são feitos os sistemas de arquivos para memórias flash. O quinto capítulo descreve o sistema de arquivos proposto (RIFFS), sua arquitetura e projeto. No sexto capítulo é mostrado como foi feita a implementação do primeiro protótipo. Por último, é mostrado uma breve conclusão do trabalho.

Capítulo 2

Memórias Flash

Este capítulo descreve uma visão geral da **Tecnologia de Memórias Flash**, trazendo principalmente o estado-da-arte neste campo e servindo como base para as próximas sessões neste texto.

A memória *flash* é um tipo de memória não volátil, mas com seu funcionamento bem distinto. Ela trabalha como a união das características de leitura e escrita das **Memórias de Acesso Randômico (RAM)** com as características de armazenamento das **Unidades de Disco Magnético**. O armazenamento dos dados dentro dessas memórias é dado em células, como nas RAMs Dinâmicas (DRAM), mas ainda trabalha como um disco magnético pelo fato da persistência de dados quando a energia é desligada. Por causa da sua alta velocidade, sua grande resistência contra impactos, seu tamanho reduzido, e seu baixo consumo de potência, as flashes tornaram-se um meio ideal para armazenamento em várias aplicações embutidas como câmeras digitais, telefones celulares, impressoras, roteadores, tocadores de MP3, etc [GRO 96].

2.1 Conceitos Gerais

Memórias flash são similares às tão conhecidas memórias EEPROM, e a principal diferença entre elas está no fato de que as memórias flash são apagadas somente em blocos, e não “por inteiro” como era feito nas EEPROMs, possibilitando assim, criar

um sistema para gerenciar dados dentro dela, uma vez que não é preciso perder toda sua informação. Pelo fato do apagamento ser baseado em setores, os circuitos das flashes acabam sendo mais simplificados, permitindo assim uma maior densidade em relação a uma memória EEPROM equivalente.

Existem atualmente vários tipos de tecnologias de memórias flash, as quais podemos citar: *NOR*, *DINOR*, *T-Poly*, *AND*, *NAND*; e cada uma dessas tecnologias requer um gerenciamento (funções de leitura, escrita e apagamento) específico [ST 01]. Estas tecnologias permitem às memórias flash, reterem dados sem uma fonte de energia por períodos longos como 20 anos, por exemplo. No entanto, essa mesma tecnologia é responsável por um dos grandes problemas das memórias flash: a **limitação no número de apagamentos**, por causa do desgaste das células de armazenamento. Com isso, os fabricantes precisam fixar o número de apagamentos que garanta a integridade dos dados (por exemplo: 100.000 apagamentos), e atualmente, este valor é razoável para a maioria das aplicações embutidas.

Pelo fato das memórias flash serem apagadas por setor, estes requerem uma atenção especial quanto ao seu tamanho. Os fabricantes, além de produzirem memórias com **setores de diferentes tamanhos**, ainda produzem *chips* com diferentes tecnologias, que também são conhecidos no mercado como **memórias flash híbridas**. Como característica, ainda podemos citar a proteção por hardware em alguns setores da flash, para melhor proteção dos dados.

2.2 Operações

Existem três operações básicas que podem ser realizadas em uma flash: **leitura**, **escrita** e **apagamento**. O tempo de leitura e escrita de uma flash é normalmente equivalente às mesmas operações em uma DRAM, mas o tempo de operação de apagamento é bem mais lenta (chegando perto da casa dos segundos). Esta limitação pode ser contornada enviando um comando de “*pausa*” que permite parar momentaneamente a execução de uma rotina de apagamento para fazer uma outra operação, voltando depois ao estado anterior.

Leitura: A leitura de uma memória flash é bastante parecida com as leituras em memórias voláteis convencionais. Para ler um dado, é só escrever o endereço desejado no “barramento de endereços” da memória, e capturar o dado do “barramento de dados”. Isso torna o acesso dos dados em memórias flash mais rápido que os meios magnéticos.

Com o intuito de reduzir os acessos e aumentar a quantidade de informação lida, alguns fabricantes implementam em suas flashes, sofisticados métodos de acesso, como: **buffer de páginas**, e **leitura seqüencial**. No primeiro método o *chip* contém uma memória volátil interna que armazena temporariamente os dados, permitindo a leitura de uma página inteira. O segundo método é conhecido como “rajada” (*burst*), onde é preciso informar apenas o primeiro endereço de uma leitura seqüencial de dados.

Escrita: O formato da escrita dos dados em uma flash é um pouco diferente daquele que estamos acostumados a pensar. Uma flash é dita “apagada” quando todos seus bits possuem o nível lógico “1” (um), e dentro desta filosofia, para escrevermos algum dado em uma flash, é necessário “trocar” alguns bits para “0” (zero). Nesse modo, podemos resumir a operação de escrita como: “escrever zeros”. É importante lembrar que o contrário: “transformar zeros em uns”, só é possível no setor inteiro ¹. Assim como na leitura, ainda temos métodos sofisticados para escrita, como: **buffer de páginas** e **escrita seqüencial**.

Apagamento: Como dito anteriormente, o apagamento se dá em um setor inteiro da flash. Esse apagamento é realizado colocando todos os bits para o valor “1” (um). Pelo fato do apagamento ser a operação mais demorada de uma flash, existe um método de **espera** para melhorar o desempenho do sistema. Com este método é possível parar a operação de apagamento para realizar outro tipo de acesso no dispositivo.

¹Se tentarmos apagar somente um dado, provavelmente a memória flash sinalizará um evento em um de seus pinos. Este evento pode ser reportado ou não, dependendo da implementação do fabricante.

2.3 Tecnologias

Desde a invenção das memórias flash, fabricantes têm procurado alternativas para aumentar o desempenho e a capacidade desses dispositivos. Como essas memórias ganharam novos mercados, tecnologias foram incorporadas para fazer das flashes um produto mais competitivo no mercado de armazenamento. O avanço destas memórias e suas tecnologias mais importantes estão citados abaixo:

Tamanho de setor variável: alguns modelos possuem setores de tamanhos variados que permitem uma melhor manipulação do sistema de arquivos. Esses setores de tamanho variados são necessários quando se deseja bloquear dados na flash. Normalmente os setores de tamanho variado em uma flash estão no início ou no fim do seu espaço de endereçamento.

Paralelismo de operações: as flashes podem ser formadas por diferentes bancos de armazenamento que operam em paralelo, e assim, operações podem ocorrer simultaneamente na mesma flash, desde que estejam sendo feitas em diferentes bancos.

Interface padronizada (CFI): “*Common Flash Interface*” (CFI) é um conjunto de operações adotado por um grupo de fabricantes com a intenção de padronizar o acesso às informações das memórias flash. Por exemplo, uma memória flash possui informações como seu número de série, número do fabricante, número de setores, etc. O CFI foi o padrão adotado para a requisição dessas informações, para que as aplicações não se preocupem em conhecer os detalhes dos dispositivos e das versões de cada fabricante.

Tecnologia de vários-níveis: essa tecnologia se refere à capacidade de armazenar dois bits de informação em apenas uma célula. Normalmente uma célula consegue armazenar apenas um bit de informação, e com esta tecnologia, o tamanho do *chip* está sendo reduzido pela metade.

Segurança: alguns modelos possuem registradores de segurança que indicam qual setor na flash pode ser “protegido”. Este termo: “segurança”, não está relacionado com

criptografia, mas sim com o bloqueamento físico de um setor na flash.

Setores Híbridos: alguns modelos de memórias flash podem ter setores de diferentes tecnologias. Estes dispositivos são chamadas de **híbridos** e encontrados em fabricantes que disponibilizam memórias com **setor de boot**. Normalmente este setor de *boot* é uma tecnologia diferente dos outros setores pelo fato do seu acesso ser mais lento.

2.4 Estudos de Casos

Esta sessão mostra alguns fabricantes de memórias flash bem como sua tecnologia disponível no mercado.

2.4.0.1 AMD

AMD disponibiliza vários modelos de memórias flash compatíveis com a interface CFI. As voltagens de seus dispositivos variam entre 1.8V até 5.0V. A densidade máxima hoje em dia é de 256 Mb com a presença, ou não, de um setor dedicado para *boot*. Suas memórias operam em temperaturas entre a faixa comercial (0 — +145° C) até super-estendida (-55 — +145° C).

Os modelos disponíveis atualmente possuem as tecnologias: “*Mirror-Bit*” e “*Dual Operation*”.

“Dual Operation”: Pela incorporação de uma memória *SRAM* no *chip*, a AMD fez um modelo que pode funcionar com várias operações simultâneas. Todas operações são executadas na memória *SRAM* e depois transferidas para flash. A aplicação não tem essa visibilidade de operação, e usa o componente como se fosse uma flash normal.

“MirrorBit Technology”: Esta é a tecnologia de **vários-níveis**, que diferencia esse tipo de flash das demais, pela implementação de uma célula que armazena dois bits de informação. Dessa maneira o *chip* tem metade do tamanho se comparado com outro de mesma capacidade.

2.4.0.2 Intel

Os produtos da Intel também possuem suporte à interface CFI. Os componentes operam a uma faixa de voltagem entre 1.8V até 5.V, e alguns modelos necessitam uma voltagem maior (aprox. 12V) para realizar operações especiais como bloquear um setor, etc. Os tipos de dispositivos fabricados podem ter as seguintes tecnologias: paralelismo (vários bancos), tamanho de setores variados, setores híbridos e células de vários-níveis. A Intel tem ainda duas tecnologias exclusivas que minimizam o tempo de escrita na flash, chamados de “*Enhanced Factory Programming*” (EFP) e “*Buffered Enhanced Factory Programming*” (BEFP).

“EFP” e “BEFP”: Essas duas tecnologias são usadas para diminuir o tempo de escrita na flash, usadas normalmente em sistemas embutidos que não possuem interação com o mundo externo após entrarem em funcionamento. Para isso é preciso pré-gravar o *chip* antes de colocá-lo no sistema definitivo. Apesar do nome diferente, essas duas tecnologias realizam a mesma tarefa, mas a diferença entre elas está no fato que o “BEFP” possui células de vários-níveis.

2.4.0.3 ATMEL

As memórias da ATMEL funcionam com uma voltagem variando entre 2.7V e 5.0V. A capacidade de armazenamento das memórias desta companhia variam desde 32 Mbits até 512 Mbits. As frequências de leitura podem chegar até 100MHz. Em alguns modelos, é possível definir o tamanho do barramento (variando entre 16 e 32 bits) em tempo de execução. As tecnologias mais usadas são: paralelismo e setor de *boot*.

Três exemplos dos principais modelos de tecnologias dessa empresa são explicados abaixo. “*Fast Programming Time*” e “*Serial Flash*” são tipos de flash que utilizam técnicas comuns para melhorar a qualidade do dispositivo. Já o modelo “*Data Flash*” possui uma tecnologia mais elaborada.

“Fast Programming Time”: é uma tecnologia para flash que possui um segundo estado de voltagem para escrita de dados. Esse segundo estado é aproximadamente 12V e isso

diminui o tempo de escrita das flashes deste modelo.

“Serial Flash”: este modelo implementa o seu protocolo de comunicação via **Interface de Periféricos Seriais** (SPI). Esta interface faz com que a flash possa ser usada como uma memória substituta das “EEPROM SPI”, sem nenhuma mudança no “layout” da placa (caso a pinagem for compatível). Este *chip* opera em 20MHz, e sua densidade varia desde 512 Kbits até 4 Mbits.

“Data Flash”: Flash compatível com a Interface SPI. Pode ter alguns setores híbridos, e o seu acesso pode ser serial ou paralelo, dependendo do número de bancos do dispositivo.

2.4.0.4 MICRON

Os produtos da MICRON são compatíveis com a Interface CFI. As voltagens dos dispositivos variam entre 2.7V e 5.0V. O tamanho do barramento pode ser escolhido em tempo de operação, entre 8 ou 16 bits. As temperaturas variam entre comercial à estendida. Alguns modelos possuem a tecnologia de *boot-sector* e multi-bancos. O principal modelo é chamado de “*Sync Flash*” que possui uma interface *SDRAM* com o mundo exterior, fazendo o software enxergá-la como uma memória do tipo RAM.

“Sync Flash”: Neste modelo, uma interface *SDRAM* é implementada. Com um alto desempenho de leitura (similar à uma RAM equivalente), esse tipo de flash tornou-se uma escolha competitiva para aplicações que necessitam executar código, ao invés de simplesmente armazenar dados. Os planos para o futuro desta tecnologia é substituir as atuais memórias RAM para persistência de programas.

Capítulo 3

Sistema de Arquivos

O Sistema de Arquivos serve para dar suporte ao armazenamento de arquivos de vários tipos, como textos, desenhos, executáveis, etc. Entre outras tarefas, o sistema de arquivos deve prover para o usuário uma interface simples e fácil de usar, na manipulação de seus dados. Ele é responsável por implementar em software um recurso que não existe no hardware. O hardware oferece simplesmente um grande conjunto de bytes contíguos, e a tarefa principal do sistema de arquivos é implementar a abstração de arquivo em cima do dispositivo de armazenamento. Este capítulo trata dos conceitos relacionados aos dispositivos de armazenamento e ao gerenciamento de Arquivos e Diretórios pelo sistema de arquivos, descritos a seguir.

3.1 Dispositivo de Armazenamento

O dispositivo de armazenamento de um sistema de arquivos pode ser qualquer mídia (também chamada de memória secundária), a qual provê um meio de armazenamento em massa, e a **persistência de dados**¹. Algumas tecnologias de dispositivos de dados persistentes, utilizam o acesso a dados através de blocos físicos, como é o caso do Compact Disk (CD), Hard Disk (HD), Floppy Disk, etc. O tamanho de cada **bloco físico** pode variar de acordo com cada fabricante, e no intuito de padronizar o acesso a

¹Persistência também pode ser entendida como a retenção de dados previamente armazenados, sem fonte de alimentação.

dados, os sistemas de arquivos implementam uma estrutura chamada de **bloco lógico**.

Outra característica importante no acesso a dados é preocupação com o **gerenciamento de blocos livres**. O sistema de arquivos precisa saber quais blocos estão ocupados e quais estão “livres”, para realizar suas operações de leitura e escrita, a fim de garantir a integridade do sistema. A seguir são mostrados os conceitos de: blocos lógicos e seu gerenciamento.

3.1.1 Blocos Lógicos

O conceito de Blocos Lógicos surgiu da necessidade em homogeneizar as operações em diferentes dispositivos. Desta forma as camadas superiores podem trabalhar com blocos lógicos de qualquer tamanho, fixo ou variado, sem se preocupar com as peculiaridades específicas de cada dispositivo, implementadas pelas camadas de acesso ao hardware. Um fator relevante para um sistema é o **tamanho** do bloco lógico utilizado. Este tamanho pode ser fixo ou variado, conforme mostrado a seguir:

Blocos de tamanho fixo: Dentro deste cenário, existe um fator muito importante na escolha do tamanho do bloco, que é a **granularidade do disco**. Por um lado, um dispositivo muito grande com blocos pequenos pode ser de difícil gestão, enquanto que um dispositivo muito pequeno com blocos grandes pode apresentar uma fragmentação interna indesejada.

Blocos de tamanho variado: Um sistema com blocos de tamanho variado apresenta uma maior flexibilidade do sistema em tempo de execução. Por outro lado, o código que implementa este tipo de característica, precisa ter um cuidado maior no **controle** de seus blocos. Existe a preocupação, a cada operação, do tamanho do bloco, o que não ocorre com blocos de tamanho fixo.

3.1.2 Gerenciamento de Blocos Livres

O Gerenciamento de Blocos Livres é uma das tarefas em um sistema de arquivos que adota a padronização de blocos lógicos, mostrados anteriormente. Esta

função é de extrema importância para o sistema, pois um simples erro pode sobrescrever uma área utilizada, chegando até a invalidar um arquivo todo.

Este gerenciamento é extremamente dependente do tipo do bloco lógico de dados adotado pelo sistema (fixo ou variado), e também do tipo de mapeamento dos blocos (descrito a seguir). No entanto, temos basicamente duas técnicas para o gerenciamento de blocos livres: **mapa de bits** e **lista encadeada**, conforme mostrado a seguir.

Mapa de Bits: Dentro do dispositivo é reservado um espaço onde será inserido este mapa. Ele consiste de uma seqüência de bits, onde a posição do bit indica o número do bloco que ele representa, e seu valor indica o estado do bloco livre ou ocupado. A vantagem deste método é a sua **simplicidade** de implementação e a forte tendência em alocar blocos contíguos. Sua desvantagem vem da dificuldade em gerenciar grandes mapas, uma vez que não podem ser carregados na memória principal.

Lista Encadeada: Consiste em manter uma lista encadeada contendo todos os blocos livres do disco. Para alocar um bloco, retira-se o primeiro da lista e para liberar adiciona-o na lista. Esta lista é grande no caso de um dispositivo vazio e normalmente ela é mantida na própria mídia. Conforme a ocupação do dispositivo aumenta, esta **lista diminui** até sua extinção, provendo seu espaço inicial para o usuário (o que não ocorre no conceito anterior). Ela é bastante eficiente em operações corriqueiras de alocação e liberação, mas pode gerar blocos seqüenciais completamente dispersos.

3.2 Gerenciamento de Arquivos

A manipulação de dados nos dispositivos pode conter um conjunto de atividades difíceis e indesejadas pelos usuários, como o cálculo da sua localização, controle de alocação, etc. A fim de tornar estas atividades transparentes, a funcionalidade do sistema de arquivos é passada aos usuários através do conceito de arquivo.

Arquivo é um conjunto de dados armazenados em um dispositivo. Cada arquivo contém dados do usuário que possuem algum significado para ele ou para o sis-

tema. Normalmente os arquivos possuem um nome dado pelo usuário para que este seja identificado entre os demais arquivos dentro do sistema. Além do nome, cada arquivo pode possuir uma série de outros atributos que são úteis tanto para o usuário quanto para o sistema, e entre os mais usuais podemos citar: **Tipo do conteúdo, Tamanho, Data e hora de criação, Data e hora de alteração**, etc.

Os arquivos são “vistos” pelo sistema através de uma estrutura chamada **descriptor de arquivo** (*file descriptor*). O descriptor é um registro no qual são mantidas as informações a respeito do arquivo. Essas informações incluem: os seus **atributos**, além de outros dados que não são visíveis aos usuários, mas imprescindíveis para que o sistema implemente as operações sobre arquivos. Um exemplo destes dados é o número lógico atribuído a cada *file descriptor*, também chamado de identificador e conhecido por “**id**”.

3.2.1 Operações

O sistema de arquivos deve prover um conjunto de operações para que o usuário manipule seus arquivos. A partir das operações básicas, muitas outras podem ser implementadas e exportadas como facilidades do sistema. Um exemplo é a operação de cópia de arquivo, a qual é implementada com as operações de leitura e escrita. Diferentes sistemas de arquivos, implementam diferentes funções básicas, mas podemos citar como as mais usuais:

- Criação(‘‘create’’): Cria um arquivo sem dados, e um descriptor lhe é associado. Caso não existam descritores disponíveis no dispositivo de armazenamento, a solicitação de criação é negada.
- Remoção(‘‘remove’’): Operação que libera os recursos associados ao arquivo.
- Abertura(‘‘open’’): A fim de acessar dados contido em um arquivo, um processo deve antes abrí-lo. Nesta operação, o descriptor de arquivo é trazido para as tabelas internas do sistema, na memória principal, para o rápido acesso.

- `Fechamento(''close'')`: Esta operação indica ao sistema de arquivos, que o processo não precisará mais acessar os dados do arquivo, e a tabela interna do sistema é atualizada.
- `Posicionamento(''seek'')`: Operação que atribui um valor ao ponteiro de dados do arquivo. Este ponteiro é utilizado nas funções de leitura e escrita.
- `Leitura(''read'')`: Operação responsável por ler dados de um arquivo. O ponteiro de dados indica onde começa a leitura (o posicionamento do ponteiro é feito através da função “seek”). Nesta função é preciso indicar ainda a quantidade de dados a serem lidos, e a posição de memória que os dados serão copiados.
- `Escrita(''write'')`: Função parecida com a “leitura”, só que nesta operação os dados são escritos. O ponteiro de dados indica onde começa a escrita. É preciso indicar a quantidade de dados a serem escritos e a posição da memória que contém os dados. Esta operação também pode ser chamada de **expansão** (“*append*”), caso o ponteiro de dados esteja na última posição do arquivo.
- `Leitura Atributos(''stat'')`: Função responsável pela visualização dos atributos de um arquivo.
- `Escrita Atributos(''chmod'')`: Esta operação possui a responsabilidade de escrever atributos em um arquivo.

3.2.2 Gerenciamento dos Blocos de Arquivos

O sistema precisa se preocupar com algumas características que todo gerenciamento de arquivos genérico deve possuir. Podemos citar, entre outras, como as mais comuns:

- Criação de arquivos com grandes dados;
- Possibilidade de acesso seqüencial a arquivos;
- Possibilidade de acesso direto a arquivos;

- Possibilidade de expansão de arquivos;
- Possibilidade de alteração do conteúdo de arquivos.

Estas e outras características são possíveis de acordo com o **mapeamento** dos dados do arquivo para os blocos lógicos, e conseqüentemente da **arquitetura do descritor de arquivo**. Este mapeamento está normalmente dentro do descritor de arquivo. É através dele que é possível encontrar os dados de cada arquivo.

O mapeamento pode ser realizado de três formas básicas (e mais uma série de formas mistas):

- **alocação contígua:** É a forma mais simples para alocar espaço em um dispositivo. Cada arquivo ocupa uma seqüência contígua de blocos. No descritor de arquivo é preciso manter o endereço do bloco lógico no qual o arquivo se inicia e o tamanho. As grandes vantagens deste método são a simplicidade do mapeamento e o pouco gasto de espaço para manter a informação dos dados do arquivo. O tempo do método “seek” acaba sendo rápido, pois é implementado com um cálculo simples de offset. A desvantagem aparece quando é preciso aumentar o tamanho do arquivo. Caso não exista blocos contíguos suficientes após o fim do arquivo, todo seu conteúdo precisa ser copiado para outra região do dispositivo que acomode todos os seus dados e mais a quantidade de blocos que se deseja expandir.
- **alocação encadeada:** Este tipo de gerenciamento de blocos serve para contornar a limitação da alocação anterior. Neste cenário, cada bloco contém no seu último dado o endereço do próximo bloco e assim por diante. Deste modo o descritor de arquivo continua o mesmo, armazenando apenas o bloco inicial e o tamanho do arquivo, mas uma parte de cada bloco físico é gasto para manter um endereço para o próximo bloco. A vantagem deste método está em permitir que qualquer bloco livre possa ser alocado a qualquer arquivo, sem uma alocação contígua no dispositivo. Como desvantagem, este método não permite o acesso direto a seus dados, fazendo com que a função “seek” seja lenta, gastando muito tempo com *I/O* para ler a lista de blocos encadeados.

- **alocação indexada:** Dentro deste tipo de mapeamento, o descritor de arquivo é implementado como uma **tabela de índices** (diferentemente das duas implementações anteriores). Neste esquema, cada entrada da tabela contém o endereço de um dos blocos que formam o arquivo. Assim é possível contornar as duas desvantagens anteriores. De um lado, ele não necessita da alocação contígua de blocos, e de outro, ele não precisa ler os blocos na operação de “seek”. Uma questão importante a ser tratada neste cenário é o tamanho da tabela de índices dentro do descritor de arquivo. Este tamanho tem que ser avaliado de tal forma que, possa ser construído arquivos grandes e pequenos sem consumir muito espaço. Uma técnica muito usada neste tipo de alocação é o uso de níveis de indireção na indexação, presente nos sistemas Unix [dO 01], como é o caso do Extended File System II (EXT2) [POI 01]. Desta maneira, a tabela de índices pode ser pequena e acomodar uma grande quantidade de dados, através de índices **diretos** e **indiretos**.

3.3 Gerenciamento de Diretórios

O termo **Diretório** pode ser entendido como sendo um conjunto de arquivos ou conjunto de referências a arquivos. Eles são úteis para organizar os arquivos no sistema, e são eles que nos permitem organizar os arquivos em grupos, facilitando sua localização.

As referências a arquivos são guardadas dentro do diretório, em forma de tabela, que por sua vez pode conter qualquer informação desejada pelo engenheiro de software. Cada linha desta tabela referência um arquivo do sistema, e a esta referência é dado o nome de **entrada de diretório** ou então **entrada de arquivo** [SIL 91].

O simples fato de como esta tabela é disposta, e quais suas informações, ditam “como” a estrutura de diretório pode ser formada. Existem diversas formas de estruturar os diretórios de um sistema, entre as mais básicas, podemos citar:

- **diretório linear:** Também conhecido como “*flat*”, é a forma mais simples de estruturar o sistema de diretório de um sistema. Neste caso o sistema possui somente um

diretório, e este corresponde a uma lista de todos os arquivos existentes no dispositivo. Como desvantagem não é possível separar os diferentes arquivos, impossibilitando o usuário de organizar seus arquivos em lugares separados, ou agrupá-los conforme sua necessidade. Neste caso, todos arquivos, tanto do usuário quanto do sistema, ficam em um mesmo lugar.

- **diretório em dois níveis:** Para dar mais flexibilidade ao primeiro sistema, esta implementação disponibiliza dois níveis de diretórios. Desta maneira, o sistema possui uma lista de diretórios, e cada diretório possui uma lista de arquivos. Assim, é possível que o usuário agrupe seus arquivos em diretórios, mas não é possível a criação do terceiro nível de diretórios.
- **diretório em árvore:** É possível estender o conceito de diretórios de tal forma que os usuários também possam criar livremente os seus próprios subdiretórios. Desta forma os diretórios são implementados dentro do sistema como arquivos. Cada arquivo precisa conter um campo **tipo** que o classificará como **arquivo de usuário** ou **arquivo de sistema**. O resultado é um sistema organizado em forma de árvore, e cada usuário tem a possibilidade de organizar seus arquivos da maneira mais conveniente.
- **diretório em grafo:** Dentro deste esquema, os diretórios continuam sendo implementados como arquivos. Dentro de uma entrada de diretório é encontrado o nome do arquivo, alguns atributos e uma referência (normalmente um número) do arquivo. Assim, pode-se ter um mesmo arquivo com dois nomes diferentes e em lugares diferentes.

3.3.1 Operações

As operações básicas mais comuns que podem ser realizadas sobre os diretórios são descritas a seguir:

- Criação(`'mkdir'`): Cria um diretório vazio. Convém lembrar que um diretório é dito vazio, quando ele possui apenas as entradas ponto `“.”` (referência

ao próprio diretório), e uma entrada ponto-ponto “.” (referência o diretório pai).

- `Remoção('`rmdir`')`: Remove um diretório vazio.
- `Inserção de ítem('`link`')`: Insere um ítem em um diretório. Seus parâmetros mais comuns são o nome do arquivo e alguns de seus atributos.
- `Remoção de ítem('`unlink`')`: Remove um ítem de um diretório.

As outras operações como leitura e escrita de atributos são responsabilidade da implementação de arquivos (sessão 3.2.1), e por isso não é mostrado nesta sessão.

Capítulo 4

Sistemas Arquivos para Memórias Flash

Atualmente as memórias flash estão sendo usadas como um padrão de armazenamento de dados de sistemas embutidos em geral. No entanto, tornar um simples *chip* de memória flash em um sistema complexo de armazenamento de dados não é uma tarefa simples. Pensando em aproveitar as vantagens da flash, e tentando contornar suas limitações, pesquisadores tiveram que desenvolver e “reciclar” conceitos para construir sistemas de arquivos para essas memórias, tornando eficiente a manipulação de dados. Esses sistemas de arquivos são normalmente implementados em dois modos: alguns desenvolvidos por inteiro [WU 94] enquanto que outros são construídos dentro de uma camada de software de acesso ao dispositivo (também conhecido como *driver*), mantendo assim uma compatibilidade com as camadas superiores dos sistemas de arquivos existentes [KAW 95].

Este capítulo mostra o uso das memórias flash em sistemas embutidos através do ponto-de-vista de sistemas operacionais, incluindo “device drivers” e sistemas de arquivos.

4.1 Conceitos Gerais

Apesar das várias vantagens da memória flash, ela apresenta algumas limitações, que podem ser visualizados como desafios para os engenheiros de software: nenhum dado podem ser **reescrito**, e ao invés disso ele tem que ser apagado antes. Para isso o **setor tem que ser apagado por inteiro**, e ainda tem que se tomar cuidado com o **número de apagamentos** que é limitado. Para contornar estes métodos, vários algoritmos e conceitos foram propostos desde o começo do mercado dessas memórias, e esta sessão trata especificamente destes algoritmos.

Sistemas de arquivos tradicionais, possuem a propriedade de atualização, provido pela natureza dos seus dispositivos (como os discos magnéticos, por exemplo). Isso faz com que os dados em um setor do disco possam ser atualizados, quantas vezes forem necessário, mas isso não acontece com as memórias flash. O esquema de atualização de dados nessas memórias pode ser conseguido de duas formas: **apagamento-e-escrita** e **remapeamento** de dados, mostrados a seguir.

4.1.1 Apagamento-e-escrita

Esta estratégia, representada pela figura 4.1, onde é mostrado o estado inicial de uma flash com um dado de nome *data_1* sendo atualizado, mostrado em 4.1(a). Para realizar esta operação é preciso: 4.1(b) copiar os dados válidos de todo o setor para um setor temporário; 4.1(c) apagar o setor; 4.1(d) copiar o novo dado e os dados do setor temporário e 4.1(e) apagar o setor temporário. Pelo fato desta estratégia sempre gastar o apagamento de dois setores a cada atualização, esta técnica não é implementada pelos sistemas de gerenciamento destas memórias. Isso sem contar na quantidade de processamento e tempo gasto para sua realização.

4.1.2 Remapeamento

O esquema de **remapeamento**, mostrado na figura 4.2 consiste em gravar a atualização dos dados em lugares diferentes dos originais, necessitando de uma

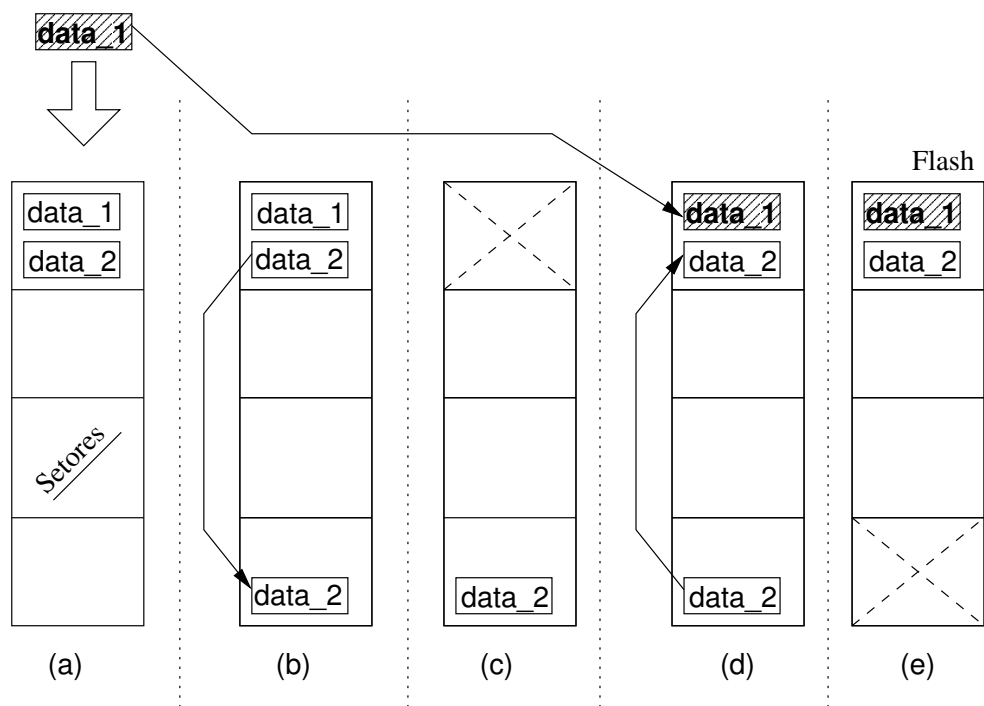


Figura 4.1: Atualização de dados na Flash.

tabela (normalmente na memória RAM) para fazer a **tradução** dos dados válidos. Na figura 4.2(a) o dado *data_1* está sendo atualizado. Dentro deste método, o dado precisa ser escrito em uma parte vazia da flash, mostrado em 4.2(b). Feito isso a **tabela de mapeamento** deve ser atualizada, e o dado antigo precisa ser **invalidado**, mostrado em 4.2(c). Como desvantagem, este método de remapeamento causa uma fragmentação nos setores por causa dos dados inválidos, necessitando assim de um **procedimento de limpeza** para apagá-los posteriormente.

Convém reforçar que a estratégia de **apagamento-e-escrita** para atualização de dados em uma flash é desaconselhada para um sistema de arquivos, porque diminui a vida útil dessas memórias, uma vez que o número de apagamento dos setores é limitado. Por esse motivo, os algoritmos de atualização de dados em memórias flash são sempre implementados através do conceito de remapeamento.

Limpeza de Setor: A estratégia de apagar um setor, reorganizando seus dados válidos em outro lugar, é chamada de **limpeza de setor**, e seu procedimento é conhecido como

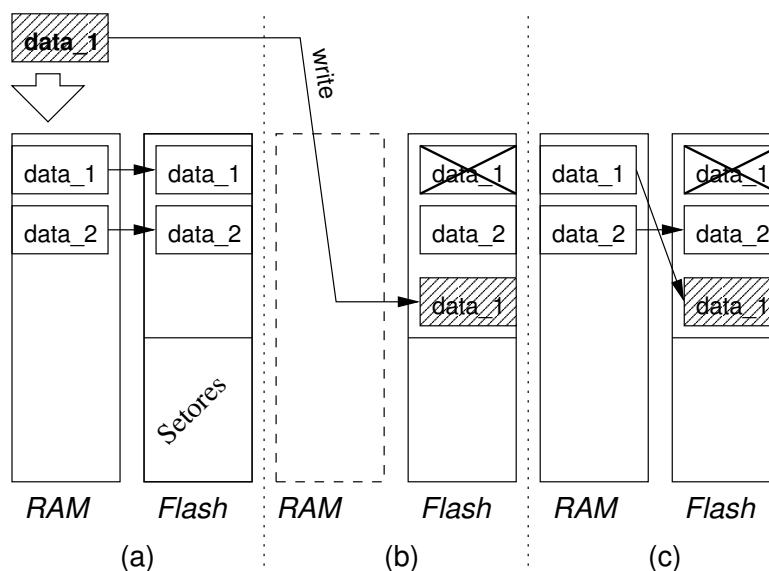


Figura 4.2: Atualização de dados na Flash.

coletor de lixo (*garbage collect*). Para realizar esta função com eficiência, muitos estudos sobre **políticas de limpeza** foram implementados. De acordo com Chiang [CHI 97], a escolha dessas políticas de limpeza possui um grande impacto na performance do sistema, podendo reduzir a eficiência de uma aplicação em até 50%.

As políticas de limpeza levam em consideração três aspectos fundamentais, como: **seleção** do segmento a ser limpo, a **reorganização** dos dados, e o **início da rotina** de limpeza. O primeiro conceito leva em conta “*quais*” e “*quantos*” segmentos devem ser limpos, com isso o coletor de lixo tem a oportunidade de trabalhar com uma maior variedade de arquivos, podendo realizar com mais eficiência a estratégia de reorganização dos dados. A reorganização se preocupa em “*como*” agrupar os diferentes dados (por exemplo, dados do mesmo arquivo ou do mesmo diretório). Já o terceiro conceito se preocupa “*quando*” será o início da rotina, que por sua vez pode ser realizada de três modos: por **tempo determinado**, por **porcentagem** de utilização da flash, ou por um **processo** de baixa prioridade no sistema que está sempre fazendo esse serviço.

Em cima desses conceitos, pesquisadores desenvolveram conceitos para ter um melhor desempenho do gerenciamento de dados dentro de uma flash. A maioria dos estudos tem se voltado para os problemas do coletor de lixo, com o objetivo de **redu-**

zir o número de apagamentos e o número de páginas copiadas.

Estratégias do Coletor de Lixo: Existem atualmente, vários esquemas para realizar com eficiência a limpeza de um setor. O primeiro método proposto foi através da política “**Greedy**”, que por sua vez recicla o setor que possui o maior número de dados inválidos. Estudos como o de Kawaguchi [KAW 95] mostram a ineficiência desse método. Esses pesquisadores propuseram então uma outra estratégia para limpeza do setor, chamada de: “**Cost-benefit**”. Dentro deste esquema, é atribuído um peso¹ (ou um valor) para cada dado escrito no setor, e o coletor de lixo executa seu método de limpeza com base nesses valores.

Chiang [MLC 99], melhorou o desempenho do coletor de lixo adotando uma política de dados “**hot-cold**”. Dentro deste esquema, os dados são agrupados de acordo com a sua taxa de atualização, sendo assim, os dados mais antigos, tendem em ficar no mesmo setor. Douglis [MAR 94], fez um estudo **estatístico** do coletor de lixo. Ele afirma que a eficiência de um sistema diminui significativamente quando a utilização da memória flash é alto. Como exemplo, ele explica que quando a utilização da flash for aumentada de 40% para 95%, o tempo de resposta das operações de escrita podem cair até 30%, e o tempo de vida da flash pode ser reduzido a um terço.

4.2 Device Drivers

Algumas memórias flash possuem um encapsulamento², exportando assim uma visão de disco magnético para o sistema operacional. Isso é necessário quando existem aplicações que precisam se comunicar com a flash, mas não se possui tempo para desenvolver uma pesquisa sobre essas memórias, e desta maneira os sistemas de arquivos para discos não precisam ser alterados. Este não é o nosso caso, uma vez que esse estudo se volta em para um eficiente gerenciamento da memória flash, via software. Com isso esta sessão mostra a primeira camada de software (presente em uma memória volátil)

¹Este peso pode ser entendido como a data de escrita dos dados

²Encapsulamento pode ser entendido como um hardware adicional.

sobre essas memórias, normalmente chamado de “*device driver*”.

4.2.1 Estudo de Casos

Os “*device drivers*” podem ser implementados de dois modos: gerenciando a flash por inteiro, exportando para as camadas superiores, uma emulação de um disco, ou exportando apenas funções básicas de leitura, escrita e apagamento, fazendo com que o sistema de arquivos seja responsável pelo gerenciamento de dados nessas memórias. Como exemplo de um driver que emula um disco, podemos citar *Flash Translation Layer*, e como exemplo da segunda estratégia podemos citar *Memory Technology Driver*.

Flash Translation Layer (FTL): FTL é um driver de gerenciamento para flashes que implementa um mapeamento de endereços lógicos para *pequenos endereços físicos*, através de uma tabela construída na memória RAM no início do sistema, de uma maneira transparente, provendo assim para as camadas superiores, a visão de um disco magnético (com blocos de aproximadamente 512 Bytes). Pela teoria, a FTL habilita qualquer sistema de arquivos a ser instalado sobre uma flash, mas normalmente quem faz muito uso desta camada são os sistemas para Windows como VFAT por exemplo.

Como mostrado na figura 4.3, a FTL divide a flash em uma ou mais **Unidades de Apagamento (UA)**, tantas quantas for o número de setores. O tamanho de uma UA depende do tamanho do setor a qual ela faz parte. Cada unidade de apagamento pode ser dividida em três partes distintas: **Cabeçalho**, **Mapa de Alocação de Blocos (MAB)** e vários **Blocos de Leitura/Escrita (BLE)**. O cabeçalho possui informação sobre a unidade de apagamento, como seu tamanho e tamanho dos BLE. O mapa de alocação se situa após o cabeçalho, e contém informações de estado e de localização de cada bloco de leitura/escrita.

Blocos de leitura/escrita, por sua vez, podem ser classificados em três tipos de dados: **Dado de Bloco Virtual (DBV)**, **Mapa de Bloco Virtual (MBV)** e **Página de Remapeamento (PR)**. Um DBV contém informações para realizar a tradução de

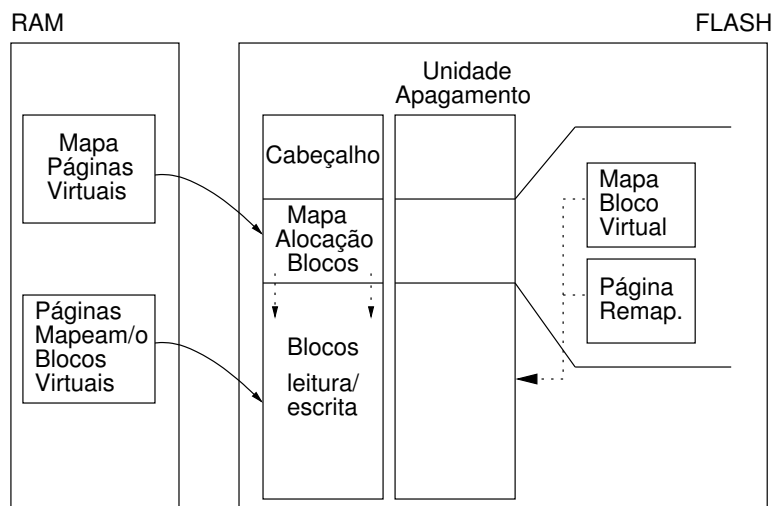


Figura 4.3: Estrutura da FTL.

endereços. Ele é organizado como uma tabela com várias entradas, sendo que cada uma aponta para um endereço físico da flash, onde os dados relacionados àquele bloco residem. O número do bloco virtual (informado pelo sistema de arquivos) é usado como índice nesta tabela. As páginas de remapeamento possuem as atualizações recentes de uma MBV, aumentando a vida útil dos setores. Alguns setores da flash são deixados como temporários (também chamados de unidades de transferência) que são usados pelo coletor de lixo para armazenar dados temporários, caso seja preciso. Estes setores são importantes quando acaba a fonte de energia quando se está realizando algum método do coletor de lixo. Com isso, garante-se que nenhum dado seja perdido durante sua execução.

Durante a inicialização do sistema, o driver FTL percorre todos cabeçalhos e todos os Mapas de Blocos Virtuais para montar na memória RAM uma tabela que reflete a estrutura do sistema. Não é necessário deixar todas MBV armazenadas na flash, e para isso o FTL disponibiliza a flexibilidade do usuário escolher a porcentagem de MBV carregada na RAM. Sendo assim, a memória RAM conteria uma estrutura chamada de *Mapa de Páginas Virtuais* (MPV) que conteria apontadores para as MBV que foram deixadas na RAM. O mesmo acontece com as páginas de remapeamento.

Memory Technology Driver (MTD): O driver MTD [WOO 02] é uma especificação de software recente dentro de projetos do sistema operacional Linux, principalmente para a área de sistemas embutidos. O projeto MTD define uma interface genérica para acesso a dispositivos de memória, em particular, dispositivos flash. Além de exportar uma interface de **pequenos blocos** para uma possível emulação de disco, o driver MTD ainda exporta uma interface de acesso a **simples caracteres**, que permite aos sistemas de arquivos possuir uma visão da flash como uma memória linear de dados³.

O foco do projeto MTD é definir uma **interface padrão** entre os dispositivos e as camadas superiores do sistema operacional. Neste sentido, alguns “*device drivers*”, dentro do projeto, são implementados apenas com funções básicas de acesso ao “*hardware*”, sem se preocupar com algoritmos de gerenciamento como um coletor de lixo por exemplo. Por outro lado, existem device drivers com um alto nível de conhecimento e gerenciamento do hardware, mas estes não fazem parte do grupo de componentes para o gerenciamento de flashes. O sistema MTD pode ser dividido em dois modos de operação: **modo “usuário”** e **modo “dispositivo”**. O primeiro se caracteriza por um conjunto de módulos que provê uma interface de alto-nível para as camadas superiores, já o segundo é um conjunto de módulos com funções simples, de acesso a dispositivos como leitura, escrita e apagamento de dados. A arquitetura do MTD é mostrada na figura 4.4.

4.3 Sistemas de Arquivos

Programas de aplicação podem possuir a funcionalidade de armazenar e buscar qualquer dado de uma memória flash através de serviços do “device driver”. No entanto, pode-se tornar inadequado o fato de, em sistemas embutidos, controlar diretamente os dados armazenados, principalmente em dispositivos que requerem uma atenção especial, como é o caso das flashes. Se diferentes aplicações precisam manipular dados aleatoriamente, ou se a manipulação de dados for muito intensa, então a instalação de um sistema de arquivos é a melhor solução.

³Esta visão de “*simples caracteres*” é usada pelo “*Journalling Flash File System*” (JFFS) mostrado na sessão 4.3.

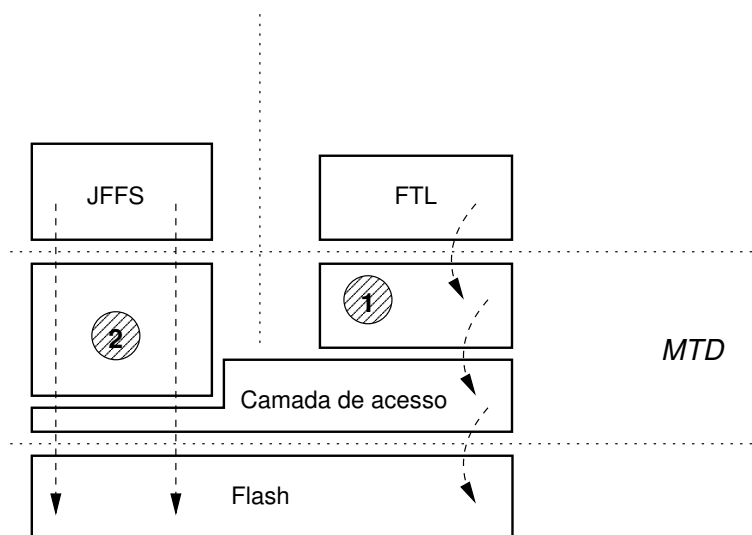


Figura 4.4: Arquitetura do MTD.

Mesmo os sistemas de arquivos sendo tão importantes no gerenciamento de dados, poucos foram propostos para memórias flash. A maioria deles, se baseiam no fato de ter uma camada que emula um disco magnético, permitindo que sistemas de arquivos não específicos para flashes sejam instalados sobre estas memórias. Felizmente, um sistema de arquivos específico possui mais vantagens, em questão de desempenho, em relação à outros sistemas classificados como “**genéricos**”, uma vez que os primeiros possuem a oportunidade de manipular diretamente as limitações impostas pela tecnologia. Esta sessão traz um estudo de casos sobre alguns sistemas de arquivos existentes atualmente.

4.3.1 Estudo de Casos

Os Sistemas de Arquivos para Memórias Flash são normalmente implementados de dois modos distintos: alguns desenvolvidos por inteiro, enquanto que outros fazem apenas a parte de gerenciamento de dados. O primeiro, implementa todos os algoritmos de gerenciamento de dados na flash e ainda exporta suas funcionalidades para as aplicações. Esses sistemas são chamados de **específicos** como dito anteriormente, e podemos citar como exemplo: o *Journalling Flash File System* (JFFS) para o Sistema

Operacional Linux, e o *Embedded File System* (Efsys) para o QNX. Como exemplo do segundo modo de implementação, temos o *True Flash File System* (TrueFFS).

Journalling Flash File System (JFFS): O *Journalling Flash File System* [DWRH 98] implementa um sistema de arquivos **específico** para memórias flash, levando em conta sistemas embutidos. A versão número um (JFFS1) foi implementado como um sistema de arquivos com estrutura-em-log, conservando o funcionamento e algumas estruturas descritas no *Log-Structured File System* (LFS) [ROS 92]. Por causa de suas desvantagens com o coletor de lixo, e pensando em acrescentar algumas características, a *Red Hat, Inc* deu início a construção da segunda versão (JFFS2) [WOO 01].

A primeira versão do JFFS possui dois tipos de estruturas: o (1) **nodo simples**, presente no dispositivo, e (2) o **nodo indexado** que se encontra na RAM, e é a estrutura onde cada nodo simples está associado. O cabeçalho da segunda estrutura contém alguns campos de controle, como: **número identificador** (id), ponteiro para seus **dados** e outras **estruturas de controle** relacionados a ele. Para controle dos dados válidos e inválidos, a estrutura de nodo simples possui um **número de versão**, onde o nodo com a versão maior significa um dado válido. Esta estratégia é usada para o sistema de arquivos saber quais dados estão apagados. Ainda nesta versão, existe uma restrição no tamanho máximo de um nodo na flash, e assim um arquivo muito grande possuirá vários nodos simples, mas apenas um nodo indexado.

O JFFS1 possui duas principais limitações: **referências estáticas**⁴ não é suportado, e **ineficiência** no coletor de lixo. A primeira limitação pode ser contornada adotando uma política de programação da aplicação que não faça uso destas referências. Já a segunda desvantagem não tem como ser contornada, tornando crítica a utilização deste sistema.

A segunda versão do JFFS melhorou as desvantagens do JFFS1, e ainda partiu para a portabilidade do seu código para todas as plataformas que possuem o Linux como sistema operacional, dando prioridade aos sistemas embutidos. Enquanto que a versão original possuía dois tipos de nodos, o JFFS2 adotou três tipos básicos: **nodo**

⁴Conhecido como “*hard link*”.

simples - diferente da primeira versão, **nodo diretório** e **nodo apagado**. Agora o nodo simples além de conter seus dados, ainda possui toda sua **estrutura de controle**. O nodo diretório é responsável pelas referências aos nodos simples. A estratégia usada para apagar um nodo, ou melhor, invalidar seus dados na flash, é simplesmente apagando a referência a ele, contida no nodo diretório. Por fim, o terceiro nodo chamado de “apagado” possui informações dos setores que foram corretamente apagados pelo coletor de lixo. Quando algum dado é escrito no setor, este deixa de pertencer à classe desses nodos.

A eficiência do coletor de lixo foi a segunda grande vantagem do JFFS2 sobre a primeira versão. Ele conseguiu uma melhor performance mudando o esquema da reciclagem de uma simples lista circular para um sistema de gerenciamento de blocos ponderado. Segundo esse método, o algoritmo de limpeza faz decisões de qual setor será reciclado, e com isso a segunda versão ganhou muito em eficiência. Ainda para melhorar a versão, foi adicionado **compressão de dados** que pode ser usado caso o usuário configure essa opção.

Para melhor gerenciar os seus dados, o JFFS2 implementa três listas de controle: **lista limpa**, **lista suja** e **lista livre**. A primeira possui apontadores para blocos com dados válidos enquanto que a segunda aponta para blocos inválidos. A lista livre possui somente nodos apagados que são os setores que não contém nenhum dado. Além dessas três listas, esta versão do JFFS ainda mantém na memória RAM um **mapa completo** do sistema de arquivos. Esse mapa é construído no início do sistema, após o sistema de arquivos ler toda a flash. As estruturas que formam este mapa são: **nodo cache** e **nodo referência**. Para cada nodo na flash existe um correspondente nodo referência na RAM, como mostra a figura 4.5.

Ainda na figura podemos observar que os nodos referências formam uma **lista encadeada mista**. O ponteiro **próximo_nodo** nos dá a visão de uma lista que pertence a um arquivo específico, e o ponteiro **próximo_físico** forma a lista encadeada com todas as estruturas de dados do sistema. Por fim, cada nodo cache representa um arquivo, e pode ser visualizado como a **cabeça** da lista encadeada ao qual pertence.

Como dito anteriormente, as estruturas na memória RAM são construídas na inicialização do sistema. Este início envolve uma operação de quatro passos:

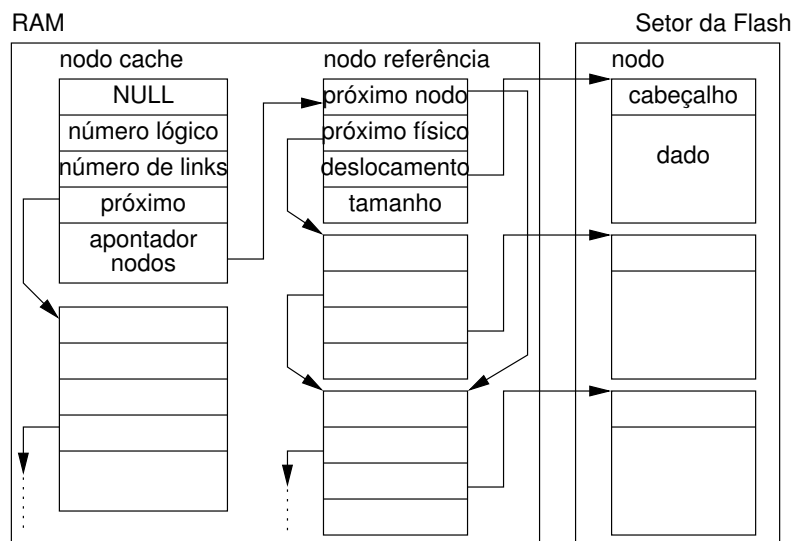


Figura 4.5: Arquitetura do JFFS2.

leitura da memória flash e alocação de todos os nodos na memória RAM, **apagamento** das estruturas que apontam para dados inválidos, **apagamento** das estruturas que não possuem referência e **apagamento** dos dados temporários. Feito isso, o sistema de arquivos começa a sua operação.

Embedded File System (Efsys): *QNX Software Systems Efsys* [QNX 02] combina as funcionalidades de um sistema de arquivos junto com um device driver. Por esse fato, existe várias versões do sistema, cada uma desenvolvida para um tipo de fabricante de memórias flash.

O software suporta dois tipos de partição: **partição simples** e **sistema de arquivos**. A primeira pode ser qualquer setor na flash que não necessite dos algoritmos de gerenciamento de dados. Como exemplo, podemos ter a imagem de um componente do QNX, que não precisará de atualização. O segundo tipo de partição, é a mais comum, onde se encontram as estruturas do sistema, seus dados de controle, etc. O formato de armazenamento de informações é proprietário, e os diretórios e arquivos são organizados como uma lista encadeada de **nodos**. Um nodo pode ser entendido como uma posição contígua de bytes em um dispositivo (pode ser uma flash, um disco, etc), e um arquivo

pode ser formado por múltiplos nodos.

Quando a flash é formatada, alguns dados de controle são escritos no setor, mas um deles fica reservado para ser usado pelo coletor de lixo como armazenamento temporário na “reciclagem” de dados. Uma característica interessante do sistema de arquivos é a sua descompactação transparente dentro da função de leitura. O mesmo não ocorre na função de escrita, onde o usuário tem que explicitamente chamar a função de compactação.

True Flash File System (TrueFFS): A empresa *M-Systems* implementou seu sistema de arquivos TrueFFS [Mic 02], baseado na camada padrão FTL, também patenteado por eles. Ele exporta a memória flash para o sistema operacional como um disco magnético. Por sua vez, não foi necessário desenvolver algoritmos de gerenciamento dessas memórias, porque o FTL provê essas funções de uma maneira transparente.

É necessário que um sistema de arquivos do tipo FAT⁵ faça a iteração entre o sistema operacional e o TrueFFS. A Figura 4.6 mostra um exemplo das camadas envolvidas nesse sistema. O TrueFFS encapsula o módulo FTL e então exporta serviços de um disco magnético, além de realizar funções específicas de acoplamento com o sistema operacional.

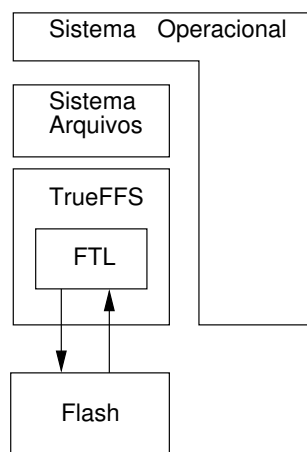


Figura 4.6: Camada TrueFFS dentro do Sistema Operacional.

⁵Encontrado em sistemas operacionais do tipo DOS da Microsoft.

Capítulo 5

Projeto do Sistema RIFFS

Através da literatura, percebeu-se que a maioria das pesquisas sobre sistemas de arquivos para memórias flash aconteceram com o intuito de aumentar o desempenho do sistema aprimorando os conceitos vistos anteriormente, como por exemplo o coletor de lixo. A proposta deste trabalho é mostrar uma *nova estrutura de armazenamento de dados em memórias de difícil atualização*, especialmente as memórias *flash*. Com isso espera-se: uma melhoria no desempenho dos “coletores de lixo” existentes, uma economia na atualização das estruturas e no processamento, e conseqüentemente um aumento na vida útil da flash.

5.1 Motivação

A seguir é mostrado alguns motivos que levaram a criação deste projeto:

Processamento de Ponteiros de Blocos: Existem sistemas que atuam dentro de drivers, como é o caso da FTL, que tentam emular pequenos setores de disco magnético (aprox. 512 bytes) sobre os setores da flash. Esta tarefa é realizada para tornar compatível as implementações anteriores de sistemas de arquivos, para que continuem acessando a memória como um dispositivo de bloco. Para cada pequeno setor emulado, é necessário um ponteiro dentro de uma tabela que na maioria dos casos é armazenada na Flash, mas mantida na RAM. Suponha como exemplo uma memória flash de tamanho

igual a 128MB. Sabe-se que cada ponteiro de bloco igual a 4 bytes, e o tamanho do bloco igual a 512 bytes. Dividindo o tamanho da flash pelo tamanho do bloco, teríamos 256K blocos. Multiplicando o total de blocos pelo tamanho do ponteiro, teríamos 512 KB sendo gastos com ponteiros em uma tabela na flash, que também pode estar replicada na RAM. A cada atualização esta tabela precisa ser alterada, e com o passar do tempo, o setor a qual ela pertence será “reciclado” pelo coletor de lixo, necessitando de uma lógica adicional para ser reescrita em outro setor. Note que não está sendo levado em conta a tabela que o sistema de arquivos cria, que também possui seu tamanho e gerenciamento.

Espaço ocupado com Ponteiros: Em uma visão mais macro de um sistema de arquivos, e ressaltando os problemas de blocos fixos, os engenheiros de software de sistemas embutidos tem que se confrontar com uma grande diversificação dos arquivos, diferente do estudo feito para sistemas UNIX, onde a maior taxa de arquivos é de tamanho pequeno [MCK 84]. Temos, cada vez mais, diferentes tipos de arquivos e tamanhos dentro de um sistema embutido. Hoje, podemos ter arquivos de alguns bytes como pequenos módulos do sistema operacional e até arquivos de alguns mega bytes como um filme dentro de uma filmadora digital. Como exemplo de desperdício com o armazenamento e gerenciamento de ponteiros para blocos lógicos (a nível de sistema), pode ser pego um pequeno filme de 100MB. Levando-se em conta que cada bloco seja de 1 KB, e cada ponteiro para bloco seja de 4 bytes, seria gasto 400KB apenas com a referência de um arquivo.

Processamento do Coletor de Lixo: Nos sistemas de arquivos que possuem a árvore de diretórios em forma de grafo, como é o caso do EXT2 [POI 01], cada diretório possui uma lista de referências para descritores de arquivo. Cada referência é chamada de entrada de diretório. Cada entrada de diretório possui algumas informações do arquivo, como por exemplo o nome e a localização do descritor de arquivo no dispositivo de armazenamento. Suponha um diretório em uma flash com com cem arquivos, por exemplo. Se o usuário por algum motivo, mudar o nome de todos os arquivos, ficariam cem referências inválidas no início dos dados do diretório, e cem referências válidas no final. Essa é a forma normal de atualização de dados dentro de uma flash, de acordo com o capítulo 4. Quando o coletor

de lixo escolher o setor no qual os dados do diretório estão inseridos, será preciso um processamento adicional para copiar apenas as referências válidas para outro setor. Com isso gasta-se tempo, memória e processamento, para concluir a operação.

Processamento de Atributos desnecessários: Outro objetivo deste projeto é a simplicidade. Quando falamos de sistemas embutidos, referimo-nos a sistemas e aparelhos, normalmente utilizados por uma pessoa. Quando acontece este caso, as informações de controle de usuário tornam-se desnecessárias para o sistema em questão. Isto também é válido para outras informações de um sistema mais complexo, como é o caso do JFFS2. Através da literatura, não foi possível encontrar meios para eliminar o controle de usuário do sistema no JFFS2 por exemplo, acrescentando assim um overhead em aplicações que não o necessitam. Muitas vezes, em sistemas embutidos, não existe a necessidade do controle de usuário, e como exemplo podemos citar uma filmadora digital que não utilize este recurso.

Para eliminar todas as características descritas anteriormente, a arquitetura do sistema de arquivos pode ser mostrada através dos modelos de *Gerenciamento de Arquivos*, *Gerenciamento de Diretórios*, e *Gerenciamento do Dispositivo de Armazenamento*, mostrados a seguir:

5.2 Objetivos

O principal objetivo deste projeto é evitar a complexidade das estruturas clássicas de sistemas de arquivos, e a atualização de dados. Para conseguir este objetivo, foi preciso uma nova arquitetura no gerenciamento de diretórios dentro das estruturas criadas na flash, chamado de árvore reversa, e o resgate do conceito de blocos lógicos de tamanho variado no gerenciamento de arquivos. Resumidamente, as características necessárias a este projeto são listadas a seguir:

- **Simplicidade das estruturas:** No início do projeto, duas das principais metas eram: tornar as estruturas físicas armazenadas na flash o mais simples possível, para economizar em espaço e evitar as atualizações, e simplificar o gerenciamento das

entradas de diretório, ao longo de sua vida. Com o começo do projeto, percebeu-se que não só o espaço estava sendo economizado, como também a vida útil dos setores.

- **Arquivos possuem todas informações:** Este requisito surgiu da necessidade em eliminar as entradas de diretório das estruturas de armazenamento, para que não ficassem resíduos de um arquivo apagado dentro de um diretório. Desta maneira, as informações da entrada de diretório foram agrupadas com o descritor de arquivo, e a esta união, deu-se o nome de *contexto de arquivo*. Convém lembrar que esta estrutura *contexto de arquivo* é encontrada armazenada na memória flash, e serve para dar suporte para a construção do sistema na memória principal.
- **Navegabilidade do sistema:** Neste projeto, de acordo com suas características, o diretório não possui uma lista de referências à descritores de arquivos. Desta forma, não existiria uma navegabilidade do sistema. A solução adotada foi acrescentar ao contexto de arquivo uma referência para o diretório ao qual ele pertence, também chamado de “diretório pai”. Pelo fato do diretório ser implementado como um arquivo dentro do sistema, ele também possui um contexto, que por sua vez aponta para seu pai, e assim sucessivamente. Assim a árvore armazenada no dispositivo acaba sendo reversa. Como não existe uma navegabilidade natural em uma árvore deste tipo, é necessário montá-la na memória principal, como uma árvore direta, garantindo então sua navegabilidade.
- **Blocos Lógicos:** Por causa da grande variação de tamanho dos arquivos em sistemas embutidos atuais, fica difícil prever qual o tamanho ideal da estrutura de dados, também chamada de *bloco lógico*. Foi pensando desta maneira que este projeto adotou uma estrutura de *blocos de tamanho variável* para o armazenamento de dados. Como gerenciamento de blocos de arquivos, o mapeamento adotado foi do tipo indexado.
- **Fragmentação:** Não existe fragmentação externa, neste projeto. Pelo fato do tamanho dos blocos ser variado, qualquer espaço pode ser alocado como um bloco

lógico de dados.

- **Portabilidade:** O projeto foi implementado através da linguagem “C++” e posteriormente portado para “C”. Como não foi feito uso de bibliotecas que vem junto com a linguagem, o código escrito pode ser compilado para qualquer plataforma (desde que exista um compilador).
- **Interface Externa:** Uma biblioteca de funções foi idealizado como produto da implementação. Desta forma, este trabalho pode ser incorporado em qualquer outro contexto, sem a necessidade de código adicional.

5.3 Modelo Arquitetural

Para realizar todas as características descritas anteriormente, a arquitetura do sistema de arquivos pode ser mostrada através dos modelos de *Gerenciamento de Arquivos*, *Gerenciamento de Diretórios*, e *Gerenciamento do Dispositivo de Armazenamento*, mostrados a seguir:

5.3.1 Sistema de Arquivos

Dentro deste projeto, o conceito de Arquivo é definido como um conjunto de dados. Estes dados podem ser tanto de controle, como de usuário. De acordo com os requisitos deste projeto, foi necessário manter todas informações a respeito do arquivo dentro dele próprio. Para conseguir esta característica, o projeto RIFFS criou uma estrutura especial chamada de “Contexto de Arquivo”. A seguir é mostrado o funcionamento do contexto, e como é a estrutura interna de um arquivo. Cada arquivo possui um tipo, que os classifica em Arquivos de Usuários e Diretórios.

5.3.1.1 Gerenciamento de Arquivos:

Dentro do gerenciamento de arquivos encontramos basicamente três estruturas: o *arquivo* (propriamente dito), o *contexto de arquivo* e os *blocos lógicos de*

dados pertencentes ao arquivo.

Arquivo: Cada arquivo, possui um Contexto, uma lista de blocos, e um tipo. O contexto é responsável por agregar informações de controle, como por exemplo o nome do arquivo. O segundo atributo guarda todos blocos de dados pertencentes ao arquivo (caso existam), e seus respectivos tamanhos. O atributo tipo classifica o arquivo perante o sistema como: *Arquivo de Usuário* e *Diretório*. Na figura 5.1(b) é apresentado um exemplo de como o arquivo é visto quando carregado na memória RAM. Já na figura 5.1(c) é possível visualizar os blocos do arquivo espalhados pela memória flash. Em 5.1(a), é mostrado a visão do arquivo na forma de um “conjunto de dados”. Nesta figura, o círculo representa um arquivo de nome “file1.txt” que possui os blocos de dados identificados no exemplo por “f_1”, “f_2”, e “f_3”.

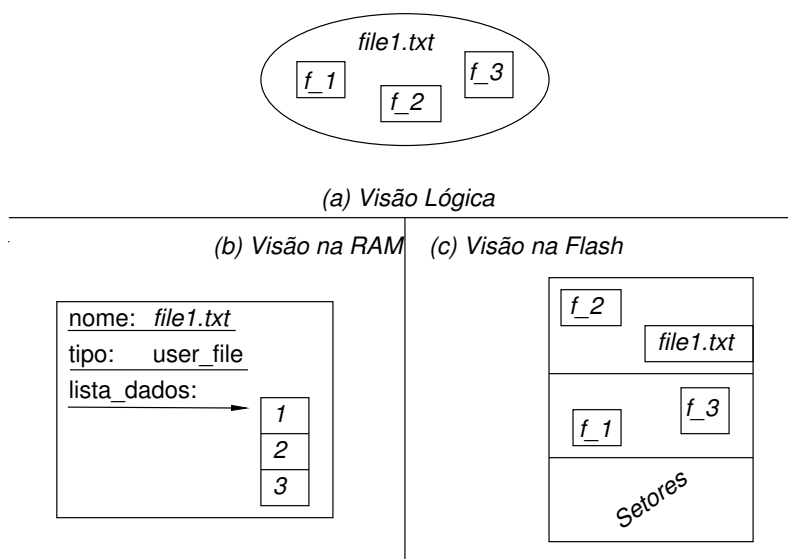


Figura 5.1: (a) Visão Lógica. (b) Visão da RAM. (c) Visão da Flash

Contexto de Arquivo: Todas informações de controle pertinentes ao arquivo estão dentro de seu contexto. O contexto é implementado como um bloco dentro da flash, e por este motivo, ele possui os mesmos atributos de um bloco lógico, descrito em 5.3.2. Ele possui em seus dados: uma referência para o “contexto pai”, e o nome do arquivo. A referência

para o “contexto pai” é utilizado no gerenciamento de diretórios, e explicado a seguir. O nome do arquivo é um array de caracteres e guarda o nome escolhido pelo usuário.

Organização dos blocos de arquivo: Cada bloco lógico de um arquivo possui uma versão que identifica as várias partes de um arquivo. Assim, ordenando a lista de blocos lógicos em uma forma crescente, temos os dados do arquivo organizados (mais detalhes sobre o campo “versão” de cada bloco lógico ver sessão 5.3.2). Isso foi necessário para permitir que os blocos do arquivo sejam escritos de uma forma aleatória em qualquer parte do dispositivo, garantindo sua reconstituição no início do sistema. A figura 5.1(a) mostra dois tipos de blocos do arquivo “file1.txt”: blocos de usuário, e um bloco do tipo contexto. Como blocos de dados do tipo usuário, temos: “f_1”, “f_2”, e “f_3”, e o bloco do tipo contexto está representado pelo nome do arquivo: “file1.txt”.

5.3.2 Dispositivo Armazenamento

Como meio de armazenamento foi utilizado uma memória flash da AMD com 2MB de tamanho. Este dispositivo não possui o conceito de blocos físicos, uma vez que manipula bytes em qualquer posição da memória. Isto é uma vantagem pois seus dados são acessados de uma forma aleatória, sempre com um mesmo tempo (pré-definido pelo fabricante), o que não ocorre nos discos (por causa de sua natureza).

Arquitetura Interna: O Setor é uma unidade muito importante no projeto de um sistema de arquivos para memórias flash, e por isso ele precisa ser analisado com cautela. A arquitetura física do dispositivo foi focado dentro do setor, e não dentro flash. Neste trabalho ele possui uma estrutura mostrada através da figura 5.2. Nela é possível visualizar três estruturas básicas: *estrutura de dados* (presente na área de dados), *estrutura de controle* (presente na área de controle), e o *cabeçalho* (presente no início de cada setor).

A primeira estrutura pode ser entendida como os próprios dados do arquivo, e representam os blocos lógico de tamanho variado. Eles são gravados no sentido do início para o fim do setor, e representado na figura 5.2 como a parte superior do desenho. Estas estruturas recebem o nome de “Raw_Data”. Como previsto, o tamanho

desses dados pode variar tanto quando se deseje, mas desde que não ultrapasse o valor máximo do setor. Caso isso ocorra, o sistema de arquivos se encarrega em gravar o restante dos dados em um outro lugar da flash. A segunda estrutura se refere ao controle dos “Raw_Data”, e são gravados no sentido do fim para o início do setor, e representado na figura 5.2 como a parte inferior do desenho. Estes possuem uma estrutura fixa, para que o método de leitura inicial do sistema seja rápido e eficiente. A terceira estrutura, o cabeçalho, possui informações de controle do sistema como por exemplo, o número de apagamentos do setor, etc. A seguir é mostrado em detalhes cada uma destas estruturas:

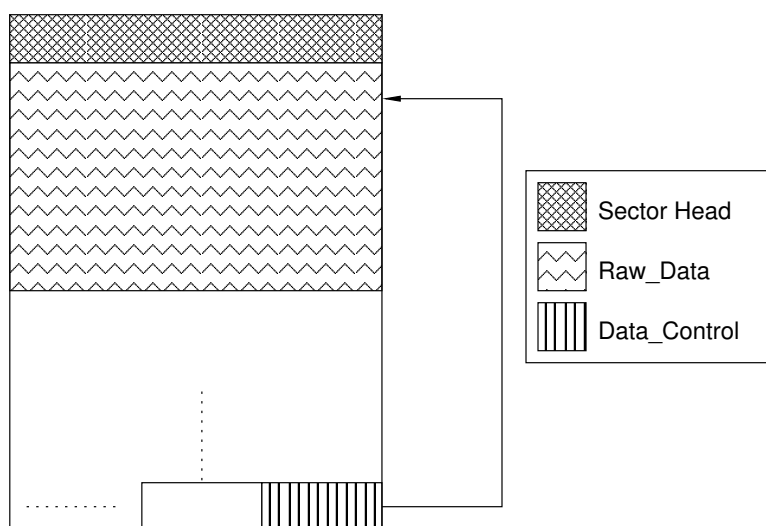


Figura 5.2: Arquitetura do Setor.

- **Área de Controle (Data_Control):** são estruturas usadas para melhorar a performance da rotina de início do sistema. Eles estão sempre presentes no final dos setores, e são gravados no sentido do final para o início. Estas estruturas contém referências aos Raw_Data dentro do setor. Os campos desta estrutura são: tamanho, offset, versão, identificador e tipo. O campo tamanho indica a quantidade de bytes que o Raw_Data ocupa. O campo identificador é o número lógico do arquivo, também chamado de file_id, ao qual o Raw_Data pertence. Este número é único dentro do sistema de arquivos. O campo offset é o local dentro do setor onde o Raw_Data está armazenado. A versão é um número correspondente à que parte

dentro do arquivo o `Raw_Data` pertence. Este campo é necessário, para distinguir as diversas partes de um arquivo, caso existam. Caso o arquivo esteja fisicamente em diferentes partes da flash, estas partes terão o mesmo `file_id`, mas versões diferentes. O campo `tipo` classifica o `Raw_Data` ao qual esta estrutura representa. Ele pode classificar o `Raw_Data` em três tipos: `User_Data`, `Log_Context` e `Context`, como pode ser visto no próximo item.

- **Área de Dados (`Raw_Data`):** é a estrutura que contém seus dados de acordo com o campo “tipo” do `Data_Control` (ao qual ele pertence). Os tipos podem ser: `User_Data`, `Log_Context` e `Context`. O `User_Data` é o dado propriamente dito, gravado do usuário. O `Log_Context` é usado caso aconteça alguma atualização com o dado. O `Context` pode ser relacionado à união de uma entrada de diretório com um `file descriptor`. Esta última estrutura possui os campos: nome do arquivo, chamado de `File_Name`, e um identificador lógico do ramo superior, chamado de `Father_Id`. Este identificador do pai é o responsável pela característica da árvore reversa. Convém ressaltar que esta estrutura `Raw_Data` corresponde ao bloco lógico de dados de tamanho variado.
- **Cabeçalho (`Sector_Head`):** estrutura presente no começo de cada setor da flash, e possui um “magic number”, um identificador (`sector_id`), e o número de apagamentos do setor (`erased_no`). O `sector_id` é um número lógico atribuído ao setor na formatação da memória, e é responsável por diferenciar os setores de um sistema de arquivos operando com mais de uma memória flash. O número de apagamentos é usado pelo coletor de lixo e pelo método de alocação, a fim de que todos setores se deteriorem por igual. Convém ressaltar que um setor é dito “vazio” quando este contém apenas o cabeçalho.

Estas estruturas são manipuladas de acordo com os módulos de software, mostrados no capítulo 6.

Blocos Lógicos: Foi adotado para este trabalho o conceito de blocos lógicos de tamanho variado. Como o conceito de blocos físicos inexistente sobre as flashes, a implementação

de blocos de tamanho variado acaba ficando mais simples. O tamanho do bloco lógico está limitado ao tamanho do setor em que este está inserido, e no caso específico da flash utilizada neste trabalho, o tamanho é de 64Kbytes. O nome dado ao bloco lógico de dados neste trabalho foi o Raw Data.

Gerenciamento de Blocos Livres: Pelo fato das flashes trabalharem com o apagamento por setor, o gerenciamento de blocos livres também segue esta filosofia. É mantido na memória principal uma lista de setores “vazios”, sendo que o primeiro ítem desta lista corresponde ao setor com o menor número de apagamentos. Para garantir que os setores se deteriorem por igual, esta lista é ordenada pelo número de apagamentos em forma crescente. Quando uma requisição chega, é alocado a quantidade de espaço requerido dentro do primeiro setor. Se caso não houver espaço suficiente, é retornado o número de dados alocados, e então é necessário que seja requisitado mais dados. Detalhes de implementação do gerenciamento de blocos livres pode ser encontrado na sessão 6.1.3.

5.3.3 Gerenciamento de Diretórios:

A arquitetura de árvore diretamente encadeada é proposto pelos sistemas de arquivos convencionais, para uma ligação das estruturas gravadas na flash. Dentro desta arquitetura, os diretórios contém informações dos arquivos e subdiretórios pertencentes a ele (caso existam). Por exemplo, no *Journalling Flash File System (JFFS2)* [DWRH 98, WOO 01], este conceito é tão utilizado, que para apagar um arquivo, basta invalidar a sua referência no diretório. Isso também ocorre em sistemas de arquivos convencionais, como o “Unix Fast File System” (UFFS) [MCK 84]. Seguindo este método de armazenamento, os diretórios são criados com um tamanho pré-definido pelo sistema. Se por um lado não existe algum arquivo pertencente ao diretório, existe um desperdício de espaço, e por outro lado, se existem vários arquivos, e a taxa de atualização é alta, existe também um desperdício de referências, só que desta vez, inválidas.

Este projeto propõe um sistema de diretórios armazenados como uma *árvore reversamente encadeada*. As estruturas pertencentes a esta árvore possuem

todas as informações a seu respeito, eliminando referências diretas no sistema. No entanto, a navegabilidade de uma árvore reversa não é adequada para um sistema de arquivos, e por esse motivo, é preciso construí-la na memória RAM de uma forma direta.

O diretório é implementado como um arquivo, e assim possui todos seus atributos, como: um Contexto, uma lista de arquivos, e um tipo (mais detalhes sobre arquivos, ver sessão 5.3.1). Pelo fato do diretório ser implementado como um tipo especial de arquivo, ele possui todos atributos de um arquivo normal. Seu contexto é idêntico ao explicado anteriormente, com os mesmos atributos. Com isso, cada diretório possui uma referência ao seu pai, e assim sucessivamente, caracterizando uma árvore reversa.

Através da figura 5.3 é mostrado três tipos de visão dos diretórios. Na figura 5.3(b) é apresentado um exemplo de como o diretório de nome “dir1” é visto quando carregado na memória RAM. A figura 5.3(a) mostra uma visão lógica da árvore reversa de contexto. É possível observar que o contexto dos arquivos “file1.txt” e “file2.txt” possuem como contexto pai o diretório “dir1”. A figura 5.3(c) mostra o conteúdo dos três contextos apresentados organizados na memória flash, de forma aleatória.

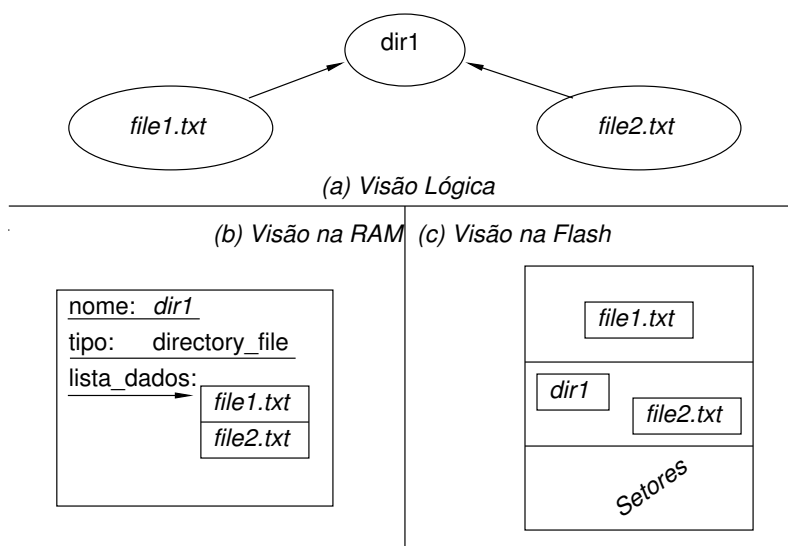


Figura 5.3: (a) Visão Lógica. (b) Visão na RAM. (c) Visão na Flash

Capítulo 6

Implementação do Sistema RIFFS

Este projeto possui sete módulos para o sistema de arquivos, conforme figura 6.1. Estes módulos são: *Flash Controler*, *Device Manager*, *Allocator*, *Garbage Collector*, *Scanner*, *File Manager*, e *Directory Manager*, explicados a seguir.

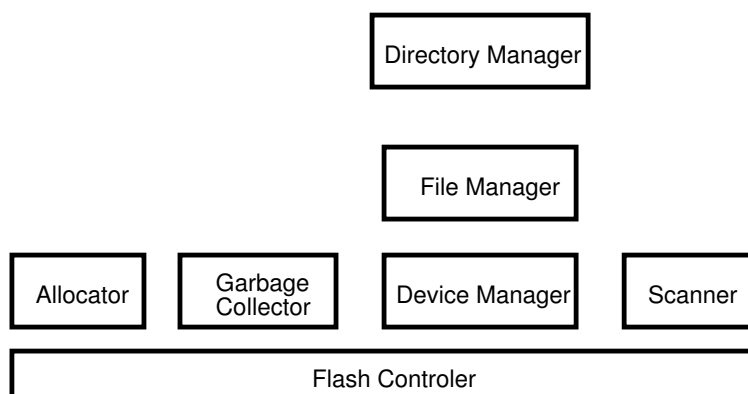


Figura 6.1: Módulos da arquitetura RIFFS.

1. **Flash Controler:** módulo responsável pela implementação das funções básicas como leitura, apagamento e escrita.
2. **Scanner:** módulo responsável pela inicialização do sistema. Percorre o final de cada setor da flash, procurando por dados válidos, armazenando-os em uma fila.
3. **Device Manager:** módulo intermediário entre árvore direta e árvore indireta. Ele é responsável pela escrita e leitura dos dados de acordo com a semântica do RIFFS.

4. **Allocator:** módulo responsável pelo gerenciamento do espaço disponível na flash. Possui uma lista de espaços vazios.
5. **Garbage Collector:** módulo responsável pela implementação das políticas de apagamento. Possui um lista de dados que precisam ser apagados.
6. **File Manager:** módulo responsável pelo gerenciamento dos arquivos propriamente dito. É subdividido em dois sub-módulos para compatibilizar as duas arquiteturas (tanto física quanto a de navegação).
7. **Directory Manager:** módulo responsável pelo gerenciamento da árvore direta. Não tem conhecimento das estruturas reversas.

6.1 Componentes do Sistema

A implementação dos sete componentes do RIFFS é mostrado a seguir.

6.1.1 Flash Control

Este módulo é responsável pela implementação das rotinas básicas de manipulação de uma memória flash. Deste jeito é possível compatibilizar todo o código de gerenciamento, com os outros módulos do sistema de arquivos para várias memórias do mercado, bastando apenas implementar esta camada.

As funções implementadas por este módulo estão descritas a seguir:

- **erase_flash:** função responsável por apagar toda a mídia.
- **erase_sector:** função responsável por apagar um setor da flash.
- **write_flash:** responsável por escrever um array de dados na flash.
- **read_flash:** função responsável por ler um array de dados da flash.

Pelo fato de diferentes fabricantes possuírem diferentes memórias no mercado, este módulo conta com a ajuda de um sub-módulo chamado de *Device_Geometry*.

6.1.1.1 Device Geometry

Este módulo tem o conhecimento da geometria da flash, que são: tamanho da flash, posição e tamanho dos setores. Estas informações podem ser obtidas de duas maneiras: implementadas estaticamente para cada modelo de flash, ou automaticamente através de uma interface padrão chamada de Common Flash Interface (CFI) [cfi 01].

Este sub-módulo implementa uma classe “Sector” que armazena informações básicas de cada setor da flash. Estas informações são: o *tamanho do setor* e seu *endereço base*. Esta classe pode ter seus atributos armazenados tanto estaticamente (via implementação) quanto dinamicamente (via interface CFI). A interface exportada por este sub-módulo utiliza esta classe como retorno de uma função. A seguir é mostrado as funções deste sub-módulo:

- **get_sector**: função responsável por retornar um “Sector”. Possui como parâmetro o índice do setor físico da flash.
- **sector_count**: retorna o número de setores.
- **flash_size**: retorna o tamanho da flash.

6.1.2 Scanner

Esta camada é utilizada na inicialização do sistema, e depois destruída. Sua principal função é construir uma lista de Data_Control que dará suporte aos demais módulos em suas inicializações.

Este módulo implementa uma classe chamada de “Data_Control”, que será mostrado a seguir, junto com as funções do módulo:

- **Data_Control**: Esta é a estrutura presente ao final de todos os setores, e contém informações sobre os Raw_Data. Os campos desta estrutura são os mesmos mostrados anteriormente (size, id, offset, type, version) no item 5.3.2.
- **init**: Função que percorre o final de todos os setores, coletando os Data_Control, e colocando em uma fila.

- **seek**: Função que posiciona o ponteiro da fila de `Data_Control`. Precisa ser passado como parâmetro o índice do ítem na fila.
- **next**: Função que retorna o próximo elemento da fila. Se for preciso acessar um ítem de índice especificado, é necessário chamar a função `Seek` antes.
- **count**: Função que retorna o número de `Data_Control` coletado durante a função `Init`.

6.1.3 Allocator

Camada responsável pelo gerenciamento do espaço livre na flash. Ela possui uma lista de setores limpos, ordenados pelo número de apagamentos, em uma ordem crescente. O algoritmo usado para implementação de alocação de espaço é o “first-fit”. Desta forma, caso haja uma solicitação de espaço, este módulo sempre pega o primeiro espaço livre dentro do primeiro setor da lista.

Este módulo exporta apenas uma função de alocação “alloc”, e implementa uma classe de controle “`Memory_Data_Control`”. Esta classe de controle é exportada pela interface do módulo e é necessária como parâmetro para a função “alloc”. Os ítems deste módulo são explicados a seguir:

- **Memory_Data_Control**: Esta classe possui esse nome porque corresponde à classe `Data_Control` (existente somente na flash), carregado na memória. Seus campos são os mesmos que o `Data_Control` (`data_offset`, `size`, `id`, `type` e `version`), e mais um campo chamado de “`control_offset`”. Este campo é usado pelo Allocator na função de alocação, e utilizado pelo `Device_Manager`, para escrita dos dados.
- **alloc**: Função responsável por implementar as políticas de alocação do sistema. Possui como parâmetro uma estrutura `Memory_Data_Control`. Seu campo “`size`” é preenchido com a quantidade de bytes que se deseja alocar. Este método se encarrega de preencher os campos “`data_offset`” e “`control_offset`”, necessários para escrever um bloco na flash (mais detalhes sobre blocos lógicos do RIFFS, ver sessão 5.3.2). Os outros campos desta estrutura são preenchidos por outros módulos. Caso

não seja possível alocar todo tamanho requerido, o campo `size` possuirá a quantidade que pode ser alocada. Caso o campo `size` seja zero, não existe possibilidade de alocação de espaço no dispositivo.

6.1.4 Device Manager

Este módulo é responsável por implementar a escrita dos dados de acordo com a semântica do RIFFS. Esta classe também pode ser entendida como a conexão entre a arquitetura das estruturas físicas (presente na flash) e a arquitetura das estruturas lógicas (presente na RAM). Este módulo exporta duas funções, uma de escrita e outra de leitura, mostradas abaixo:

- **write_data:** Esta função possui dois parâmetros: *buffer* e *Memory_Data_Control*. O *buffer* contém os dados que se deseja escrever na flash, e o *Memory_Data_Control* contém informações básicas a respeito do *buffer*.
- **read_data:** Esta função possui os mesmos parâmetros da “write_data”, mas ela realiza a operação de leitura da flash. Os dados lidos são copiados para o *buffer* da memória passado como parâmetro.

6.1.5 File Manager

Módulo responsável por agrupar as diversas partes dos arquivos (caso existam) em estruturas do tipo “File”. Após sua inicialização uma lista de arquivos é construída, que servirá para inicialização do módulo *Directory_Manager*, explicado a seguir.

O componente “File” é definido por este módulo e possui os seguintes campos: *type*, *id*, *context* e *list*. O campo “type” possui o tipo do arquivo e classifica-o em “USER_FILE” e “DIRECTORY”. O campo *id* é o identificador lógico do arquivo. O campo “context” contém informações de controle do arquivo e, por último, o campo “list” é um ponteiro para uma lista que deve ser utilizada de acordo com o campo “tipo”. Caso o tipo seja “USER_FILE”, a lista contém referências para estruturas do tipo *Raw_Data* e também é chamada de “data_list”. Caso o tipo seja “DIRECTORY”, a lista contém

referências para classes do tipo “File”, e neste caso é chamada de “file_list” (explicado em 6.1.6).

As funções implementadas por este Módulo, são mostradas abaixo:

- **init:** Função executada apenas no início do sistema. Percorre toda a lista de Data_Control do Scanner, agrupando os blocos dos arquivos. Os blocos do tipo “Raw_Data” que contém o mesmo identificador, são acrescentados na lista “data_list” do arquivo correspondente. Este procedimento é o mesmo para os blocos do tipo context, com excessão que eles são adicionados em um campo especial do arquivo: “file_context”. Depois de ter percorrido toda a lista do Scanner, é feito uma verificação dos arquivos inválidos: arquivos marcados como apagados, ou que não possuem um bloco “file_context” associado, são removidos.
- **create_file:** Esta função cria a estrutura “File” na memória e associa um identificador ao arquivo. Os atributos “data_list” e “context” são iniciados com zero e o campo “tipo” é iniciado com “USER_FILE”. O retorno da função é a posição de memória que a estrutura foi criada. Se o retorno da função for zero, significa que ocorreu algum erro.
- **remove_file:** Função responsável por gravar os atributos do arquivo como removido. O campo “versão” do bloco de dados “context” é gravado com o valor zero. Isto indica que um arquivo está apagado, pois não existe contexto com versão zero. Feito isso, uma função do módulo Garbage_Collector é chamada, pois este é o módulo responsável pelo apagamento das estruturas.
- **open_file:** Função que abre um arquivo, e aloca recursos no sistema de arquivos. Possui como parâmetro o “id” do arquivo que se deseja abrir. O retorno da função é uma classe “File” correspondente ao arquivo.
- **close_file:** Função que libera os recursos alocados pela abertura do arquivo.
- **seek:** Função responsável pelo posicionamento do ponteiro de dados do arquivo. Possui um parâmetro “position” que indica a nova posição.

- **read:** Função que lê uma quantidade de bytes de um arquivo, e armazena em uma posição de memória. O tamanho de bytes e a posição de memória são passados como parâmetro.
- **append:** Insere uma quantidade de bytes no final de um arquivo que estão em um endereço de memória.
- **truncate:** Função que retira os bytes que estão referenciados pelo intervalo com início no ponteiro de dados (posicionado pela função “seek”) até o fim do arquivo.

6.1.6 Directory manager

Módulo responsável pela construção e manipulação de uma árvore direta na memória principal do sistema. No início do sistema, este módulo percorre a lista de arquivos do File_Manager para então construir uma árvore direta na RAM. O Diretório é implementado como um arquivo do tipo “Directory”. O primeiro diretório é chamado de “raiz” ou “root” e seu nome é o caracter “/”. Este diretório não está gravado na flash. Ele possui um identificador padrão de número 1. Desta forma, todos arquivos que contém o campo “father_id” igual à 1 pertencem ao diretório raiz, e então são acrescentados à lista de arquivos do diretório. O mesmo ocorre com os subdiretórios, e assim sucessivamente. As funções implementadas por este módulo são descritas a seguir:

- **init:** Esta função é executada somente no início do sistema. A lista de arquivos do módulo “File_Manager” é percorrida, e então uma árvore direta é construída na RAM.
- **create_dir:** Esta função é similar a função “create_file”, mas com a exceção de que a estrutura de retorno possui o campo tipo com o valor “Directory”.
- **remove_dir:** Esta função remove um diretório vazio do sistema. Sua operação é idêntica a função “remove_file”.
- **link:** Esta função é responsável por inserir um arquivo ao diretório. O arquivo é associado ao diretório através de um nome, que é passado como parâmetro para a

função. Convém lembrar que cada operação de “link” sobre um arquivo sobrescreve o seu contexto anterior. Desta maneira, é implementado de uma maneira transparente as funções de renomear o nome de um arquivo e mover um arquivo para outro diretório.

- **unlink:** Esta remove a referência “father_id” do contexto do arquivo, e remove a referência do arquivo da lista de arquivos.
- **lookup:** Esta função procura por um arquivo de nome “file_name” dentro da lista de arquivos do diretório “dir”. O retorno da função é uma referência ao arquivo de nome “file_name”. Caso o arquivo não exista, é retornado um erro.

6.1.7 Garbage Collector

Pela característica da flash, quando um dado precisa ser atualizado, ele é invalidado e reescrito em outro lugar. Assim, dentro do dispositivo podem existir resíduos de dados atualizados, que precisam ser apagados posteriormente. Este módulo é o responsável pelo gerenciamento dos dados “inválidos” para o sistema. O Garbage_Collector possui duas listas: *erased_list* e *invalid_list*. A primeira lista é formada por objetos do tipo “Sector_Header”, e ordenada de forma crescente pelo campo “erased_times”. A segunda lista é formada por estruturas do tipo “Invalid_Data_Statistic” e ordenada de forma decrescente pelo campo “invalid_data_size”. A lista “erased_list” é usada quando se deseja descobrir qual o destino dos dados válidos de um setor que está sendo reciclado. A lista “invalid_list” é usada para selecionar o setor que possui a maior quantidade de dados inválidos para então ser reciclado. O primeiro elemento da lista “invalid_list” é o setor que será escolhido para ser reciclado, quando o método “execute()” for invocado. De acordo com a literatura o método de escolha de setores “origem” e “destino” usado neste trabalho é chamado de “Greedy”. A seguir é mostrado as funções e classes implementadas por este módulo:

- **Invalid_Data_Statistic:** Esta classe é responsável por armazenar uma estatística da quantidade de bytes inválidos de cada setor. Ela possui dois cam-

pos: “invalid_data_size” e “sector_id”. O primeiro campo possui o total de bytes inválidos referente ao setor identificado pelo segundo campo.

- **init:** Esta função é executada somente no início do sistema. Ela percorre toda a lista do módulo “Scanner” para detectar os dados inválidos, e começar suas estatísticas.
- **remove_data_control:** Esta função é usada quando se deseja retirar um bloco de dados aleatório do sistema. O bloco não é removido fisicamente no mesmo instante da invocação da função, mas sim quando for realizado a reciclagem do setor. É passado como parâmetro para esta função um objeto do tipo “Data_Control”, e através dos campos deste objeto, é possível coletar informações suficientes para removê-lo. Um bloco de dados físico na flash é dito “apagado” quando a estrutura “Data_Control” presente no final do setor, possuir seu campo “id” igual a zero. Desta maneira é possível diferenciar um bloco de dados “válido” de outro “inválido”.
- **remove_file:** Função responsável por retirar um arquivo do sistema. Possui como parâmetros o arquivo que se deseja remover. O arquivo em questão não é removido fisicamente no mesmo instante da chamada da função, mas sim na reciclagem do setor em que estão seus blocos de dados. A implementação desta função apaga cada um dos blocos de dados pertencentes ao arquivo através da função “remove_data_control”.
- **execute:** Função que inicia a execução do Garbage_Collector. Ela pega o primeiro setor da lista invalid_list (seleção do setor de origem), e o primeiro setor da lista erased_time (seleção do setor de destino), e inicia a cópia dos dados válidos. A cópia é feita a partir do setor de origem para o setor de destino. Para cada bloco de dados copiado de um setor para outro, é preciso atualizar a sua referencia, dentro do arquivo na memória RAM. Este conceito foi visto no gerenciamento de dados em flash, chamado de remapeamento.

6.1.8 Notas

O compromisso deste trabalho é fazer o armazenamento de árvore reversa de diretórios, sem comprometer a navegabilidade de um sistema de arquivos. Como pontos “positivos” podemos citar a baixa fragmentação causada por este projeto, pois na flash é escrito somente o necessário. Como pontos negativos este primeiro projeto não suporta “hard links”. Outra desvantagem encontrada foi no apagamento de arquivos. Nele é preciso percorrer todos os dados da lista de blocos de um arquivo, invalidando-os.

6.2 Resultados do Sistema

Após a construção do primeiro protótipo, foram realizados alguns testes práticos entre o sistema RIFFS e o JFFS2. Entre estes testes, a comparação entre a banda de escrita dos dois sistemas provê o resultado mais interessante. Esta sessão mostra alguns dados do sistema RIFFS, e algumas comparações com o sistema JFFS2.

6.2.1 Plataforma de testes

Para cada sistema de arquivos foi usado uma plataforma diferente. Esta solução encontrada foi pelo motivo de não conseguir uma versão do JFFS2 para a plataforma RIFFS e vice-versa. Ainda por este motivo foi preciso achar algum meio de compatibilizar os testes de dois sistemas em plataformas diferentes. O meio encontrado foi achar uma **porcentagem** ou uma **proporção** do tempo de escrita de dados em uma flash (através de simples funções dado pelo fabricante), pelo tempo de escrita de dados através do sistema de arquivos. Com isso obtém-se uma proporção da defasagem de tempo de escrita em arquivo em relação a simples funções de escrita na flash. A seguir é mostrado em detalhes as plataformas de testes e os cálculos realizados.

Plataforma RIFFS: A plataforma de testes utilizada neste projeto foi desenvolvida por uma empresa de Florianópolis, que atua na área de construção de sistemas embutidos para telecomunicações, chamada Khomp Solutions. O nome da plataforma é KMB01, e será

usada futuramente como uma plataforma de convergência entre o “mundo” de voz e o “mundo” de pacotes. A KMB01 possui um **processador PowerPC 405GP de 200MHz**, da IBM. No barramento do processador estão conectados: a memória flash, a memória RAM, além de outros componentes não citados por não fazerem parte do escopo deste trabalho. A memória flash possui uma conexão com o processador através de um barramento de 16 bits e seu tamanho é igual a 2MB. O tamanho da memória RAM no sistema é de 16MB. Atualmente não existe um porte de sistema operacional para esta plataforma, e por este motivo os testes do sistema RIFFS realizados foram funções simples.

Plataforma JFFS2: A plataforma utilizada para fazer os testes do sistema JFFS2 foi um iPaq H3600 da Compaq. O H3600 possui um **processador StrongArm SA-1110 de 206MHz**, da Intel. No barramento do processador estão conectadas: a memória RAM e a memória flash. A memória flash possui uma conexão com o processador através de um barramento de 32 bits, e seu tamanho é 16MB. O tamanho da memória RAM do sistema é de 32MB. Os testes realizados no JFFS2 rodaram sobre o sistema operacional Linux, e por este motivo, sabe-se a princípio que existe a presença de um escalonador e de alguns processos do kernel que não podem ser desativados. Desta maneira os testes não refletirão com **exatidão**, a comparação dos dois sistemas, mas mesmo assim, servirá como base para uma análise mais avançada.

6.2.2 Escrita de Dados

A realização destes testes baseia-se na porcentagem de desempenho perdido do sistema em relação à simples operações da flash. Este esquema possui uma fórmula bem simples explicada da seguinte maneira: suponha o Tempo de Escrita na flash realizado por Simples funções, sem intervenção de módulos do sistema, chamado de “TES”, e suponha o Tempo de Escrita da mesma quantidade de bytes através de um sistema de Arquivos, chamado de “TEA”. Dividindo o Tempo de Escrita em Arquivos (TEA) pelo Tempo Escrita Simples (TES) obtém-se a porcentagem do consumo de processamento adicional imposto pelo sistema.

Esta porcentagem de desempenho foi necessária porque os testes realizados foram feitos em duas plataformas diferentes. Os resultados dos tempos são descritos a seguir:

Tempo dos testes RIFFS: O Tempo de Escrita Simples (TES) para a flash da plataforma KMB01 foi adquirido da seguinte maneira:

- Escrever cem vezes, buffers de tamanho 1MB, e calcular seus respectivos tempos.
- Como a média dos tempos, foi encontrado: $TES = 7,307456$ seg. (para buffers de 1MB).
- Como os testes são menores que 1MB, então o resultado é dividido por 1024, para achar o tempo médio gasto para escrita de buffers de 1KB. Logo: o $TES = 7,13619$ mseg. (para buffers de 1KB).

O Tempo de Escrita em Arquivo (TEA) para o sistema RIFFS, possui seis valores. Cada valor é calculado de acordo com os testes realizados com diferentes tamanhos de buffers para escrita. Os tamanhos escolhidos foram: 16KB, 32KB, 64KB, 128KB, 256KB e 512KB. O tempo para cada tamanho de buffer foi conseguido da seguinte maneira:

- Realizar cem vezes, um “append” do tamanho do buffer em questão.
- Calcular a média do tempo gasto nestes testes.

Os tempos calculados para cada buffer, é mostrado na tabela 6.2.2, a seguir:

Tempo dos testes JFFS2: O Tempo de Escrita Simples (TES) para a flash da plataforma iPAQ H3600 foi adquirido da seguinte maneira:

- Escrita de um buffer de tamanho 12MB, e calcular seu respectivo tempo.
- Com o tempo, tempos: $TES = 53,557377$. (para buffer de 12MB).

Tamanho do Buffer	TEA RIFFS (mili-seg)
16KB	114
32KB	236
64KB	501
128KB	965
256KB	1930
512KB	3739

Tabela 6.1: Tempo de Escrita em Arquivo (TEA) para o RIFFS

- Dividindo o tempo por 12288 (12KB), encontra-se: $TES = 4,43989$ mseg. (para buffers de 1KB).

O TEA para o JFFS2 foi conseguido da mesma maneira do TEA para o RIFFS (mostrado acima), com excessão que ele foi executado na plataforma iPAQ H3600. A tabela 6.2.2 mostra os valores:

Tamanho do Buffer	TEA JFFS2 (mili-seg)
16KB	120
32KB	220
64KB	450
128KB	860
256KB	1600
512KB	2900

Tabela 6.2: Tempo de Escrita em Arquivo (TEA) para o JFFS2.

Comparação dos tempos: Após realizados os testes para aquisição dos tempos, utilizou-se a fórmula TEA/TES como cálculo de porcentagem, mostrado na tabela 6.2.2, a seguir:

Tamanho do Buffer	RIFFS - TEA/TES	JFFS2 - TEA/TES
16KB	0,46%	68,92%
32KB	3,66%	54,85%
64KB	9,70%	58,37%
128KB	5,70%	51,33%
256KB	5,67%	40,77%
512KB	2,35%	27,57%

Tabela 6.3: Comparação de desempenho entre: RIFFS e JFFS2.

6.2.3 Codificação do sistema

Atualmente, o protótipo do sistema foi construído como uma biblioteca estática, compilado inicialmente para a plataforma de testes. Esta biblioteca não utiliza nenhuma outra biblioteca externa, nem chamadas ao sistema operacional, podendo ser incorporada em qualquer outro projeto ou sistema, sem problemas de dependências.

6.2.4 Tamanho de estruturas

Abaixo é descrito o tamanho de algumas estruturas deste sistema:

- **Tamanho de informações de controle:** As Informações de Controle (IC) neste projeto são: o cabeçalho de cada setor (“Sector_Header” - SH) e a estrutura que controla o cada bloco de dados (“Data_Control” - DC). O SH está presente no começo de cada setor e possui o tamanho de 64 Bytes. O DC está presente no final de cada setor e cada unidade corresponde a um bloco gravado no setor ao qual ele ocupa, e possui o tamanho de 128 Bytes.
- **Tamanho de arquivo:** O tamanho de um arquivo dentro deste projeto é dado pelo somatório do tamanho de cada bloco de dados pertencente ao arquivo. O campo “version” da classe “Data_Control” define cada bloco distinto de dados de um arquivo. Este campo reserva 15 bits como endereço do bloco, e assim sendo o Número

Máximo de Blocos (NMB) de um arquivo pode ser:

$$NMB = 2^{15}.$$

$$NMB = 32K.(1)$$

O tamanho máximo de um bloco é limitado fisicamente pelo tamanho de um setor, ou seja, um bloco de dados não pode ser maior que um setor. Para fazer a conta do tamanho máximo de um arquivo, tomemos como exemplo uma memória flash com todos setores iguais a 64KB e tamanho suficiente para acomodar todo arquivo. Neste exemplo o Tamanho Máximo de um Bloco (TMB) pode chegar à 64KB menos Informações de Controle (IC). Logo:

$$TMB = 64KB - IC.(2)$$

Multiplicando o tamanho máximo de um bloco, pelo Número Máximo de Blocos de um arquivo, obtemos o Tamanho Máximo de um Arquivo (TMA) suportado por este projeto. Assim:

$$TMA = NMB * TMB.(3)$$

Substituindo (1) e (2) em (3), temos:

$$TMA = 32K * (64KB - IC).(4)$$

Sendo as informações de controle um “Sector_Header” e apenas um “Data_Control” por setor, temos:

$$IC = SH + DC.$$

$$IC = 64 + 128.$$

$$IC = 192.(5)$$

Substituindo (5) em (4), obtém-se:

$$TMA = 32K * (64KB - 192).$$

$$TMA = 2GB - 6MB.$$

$$TMA = 1,94GB.$$

- Número de arquivos: O número de arquivos do sistema é dado pela soma de todos os arquivos. O campo “id” da classe “Data_Control” define um arquivo distinto dentro do sistema. Neste campo são reservados 32 bits para o número lógico de cada arquivo, com exceção do número lógico zero que possui uma finalidade de controle. Assim sendo, o Número Máximo de Arquivos (NMA) que este projeto possui é:

$$NMA = 2^{32} - 1.$$

$$NMA = 4G - 1.$$

Capítulo 7

Conclusão

Este trabalho apresentou as características do projeto RIFFS (Reverse Indirect Flash File System), um sistema de arquivos para memórias flash, focado nas características destas memórias e nas estruturas de controle do software. A maior parte dos esforços se concentraram em um modo de simples estruturas escritas nas flashes.

Existem atualmente diferentes implementações para memórias flash no mercado. Algumas implementam apenas a camada de acesso ao dispositivo, como é o caso do FTL, enquanto que outros desenvolvem o sistema por inteiro, como o JFFS2, por exemplo. O primeiro tipo de sistema provê uma interface semelhante à dos discos, para compatibilização dos sistemas existentes, que não levam em conta características específicas da flash. O outro tipo de sistema é construído por inteiro e considerando as limitações das memórias. Este último tipo mistura o armazenamento de suas estruturas físicas com a navegabilidade do sistema de arquivos. Com este projeto foi possível uma validação do conceito de uma estrutura de árvore reversa para armazenamento de dados diferente dos modos convencionais.

Este projeto conseguiu provar que as estruturas escritas em memórias de difícil atualização, não precisam ser complexas, e também não precisam conter todas as informações pertinentes à operação do sistema. Com um gerenciamento de blocos de tamanho variado no gerenciamento de arquivos, e árvores reversamente encadeadas no gerenciamento de diretórios, conseguiu-se um sistema de arquivos com praticamente as

mesmas funcionalidades dos demais.

Atualmente existe o primeiro protótipo do sistema construído em forma de uma biblioteca de funções. Pelo fato desta biblioteca não possuir dependências externas, este projeto é facilmente portátil para outras plataformas (dependendo apenas do compilador). A desvantagem deste projeto, em relação aos outros sistemas de arquivos para memórias flash, é a falta de suporte para “hard-links”.

Os resultados obtidos até o momento são muito satisfatórios. Os testes sobre o “Garbage Collector” não puderam ser feitos na prática, por não ter tempo suficiente para isolar o coletor de lixo do JFFS2, mas na teoria ele funciona mais rápido, devido a simplicidade das operações. Nos teste de “escrita de arquivos” pequenos (aprox. 16KB), o RIFFS conseguiu um desempenho muito superior ao sistema JFFS2, também devido à simplicidade de suas operações.

Como trabalho futuro, existe a necessidade de integração com um sistema operacional. Isto já está acontecendo, dentro do Laboratório de Integração de Software e Hardware (LISHA) através do projeto EPOS [FRö 01] (Embedded Paralell Operating System). Desta maneira o RIFFS possuirá todas características de um sistema de arquivos, e inserido dentro do contexto de sistema operacional.

Referências Bibliográficas

- [cfi 01] Common flash memory interface specification. 1 AMD Place, Sunnyvale, CA, USA: Intel Corporation and Advanced Micro Devices and Fujitsu Limited and Sharp Corporation, 2001. Relatório técnico.
- [CHI 97] CHIANG, M.-L.; LEE, P. C. H.; CHANG, R.-C. Managing flash memory in personal communication devices. In: IEEE INTERNATIONAL SYMPOSIUM ON CONSUMER ELECTRONICS (ISCE'97), 1997. **Proceedings...** Singapore: [s.n.], 1997. p.177–182.
- [dO 01] DE OLIVEIRA, R. S.; DA SILVA CARISSIMI, A.; TOSCANI, S. S. **Sistemas Operacionais - Série Didática**. 1. ed. Sagra Luzzatto, 2001. 233 p.
- [DWRH 98] DAVID WOODHOUSE RED HAT, I. **JFFS: The Journalling Flash File System**. Red Hat, 1998.
- [FRö 01] FRÖHLICH, A. A. M. **Application-Oriented Operating Systems**. Berlin, Ge: GMD Research Series, 2001.
- [GRO 96] GROSSMAN, S. Future trends in flash memories. In: PROCEEDINGS OF THE 1996 IEEE INTERNATIONAL WORKSHOP ON MEMORY TECHNOLOGY, DESIGN AND TESTING (MTDT'96), 1996. **Proceedings...** Singapore: IEEE, 1996.
- [KAW 95] KAWAGUCHI, A.; NISHIOKA, S.; MOTODA, H. A flash-memory based file system. In: USENIX TECHNICAL CONFERENCE, 1995. **Proceedings...** New Orleans, LA: [s.n.], 1995. p.155–164.
- [MAR 94] MARSH, B.; DOUGLIS, F.; KRISHNAN, P. Flash memory file caching for mobile computers. In: IN PROCEEDINGS OF THE 27TH HAWAII INTERNATIONAL CONFERENCE ON SYSTEMS SCIENCE, 1994. **Proceedings...** : [s.n.], 1994. v.Vol. I: Architecture, p.451–460.
- [MCK 84] MCKUSICK, M. K. et al. A fast file system for UNIX. **ACM Transactions on Computer Systems (TOCS)**, [S.l.], v., 1984.

- [Mic 02] Microsoft Co. **Linear Flash Memory Devices on Microsoft Windows CE 2.1**, Online. ed., Dezembro, 2002.
- [MLC 99] MEI-LING CHIANG, PAUL C.H. LEE, R.-C. C. Cleaning policies in mobile computers using flash memory. **The Journal of Systems and Software**, [S.l.], v.48, n.3, p.213–231, Nov, 1999.
- [POI 01] POIRIER, D. **The Second Extended File System: Internal Layout**. GNU Free Documentation License and Free Software Foundation, 2001.
- [QNX 02] QNX Neutrino OS. **QNX Momentics Non-Commercial Documentation Roadmap**, Online. ed., Dezembro, 2002.
[http://www.qnx.com/developer/docs/momentics_nc_docs/roadmap/].
- [ROS 92] ROSENBLUM, M.; OUSTERHOUT, J. K. The design and implementation of a log-structured file system. **ACM Transactions on Computer Systems (TOCS)**, [S.l.], v.10, n.1, p.26–52, 1992.
- [SIL 91] SILBERSCHATZ, A.; PETERSON, J. L.; GALVIN, P. B. **Operating System Concepts**. 3. ed. Reading, USA: ACM Press, 1991. 696 p.
- [ST 01] ST Microelectronics. **Background Information: Flash Memories**, Online. ed., Fevereiro, 2001. [<http://us.st.com/stonline/press/news/back2002/b979m.htm>].
- [TEN 00] TENNENHOUSE, D. Proactive Computing. **Communications of the ACM**, [S.l.], v.43, n.5, p.43–50, Maio, 2000.
- [WOO 01] WOODHOUSE, D. JFFS: The Journalling Flash File System. In: PROCEEDINGS OF THE OTTAWA LINUX SYMPOSIUM 2001, 2001. **Proceedings...** Ottawa, Canada: [s.n.], 2001.
- [WOO 02] WOODHOUSE, D. **Memory Technology Device**, Online. ed., Dezembro, 2002.
[<http://www.linux-mtd.infradead.org/tech/>].
- [WU 94] WU, M.; ZWAENEPOEL, W. envy: a non-volatile, main memory storage system. In: PROCEEDINGS OF THE SIXTH INTERNATIONAL CONFERENCE ON ARCHITECTURAL SUPPORT FOR PROGRAMMING LANGUAGES AND OPERATING SYSTEMS, 1994. **Proceedings...** San Jose, California, United State: ACM Press, 1994. p.86–97.