

## VETORIZAÇÃO DE PROGRAMAS - A EXPERIÊNCIA DA PETROBRÁS

Simplício Lopes de Freitas  
PETROBRÁS - DEPEX - DISEP  
Av. República do Chile 65, sala 1514  
CEP 20035 - Rio de Janeiro - RJ

### RESUMO

O presente trabalho apresenta uma visão de como o Departamento de Exploração da PETROBRÁS aumentou sua capacidade computacional, de modo a processar um volume maior de dados, e implementou novos algoritmos na prospecção de óleo e gás. São também descritos trabalhos recentes para a conversão de aplicações que operam com o IBM 3090-20E/2VF e apresentados resultados de "benchmarks".

### ABSTRACT

This paper presents a view of how the PETROBRÁS Exploration Department has increased its computational capacity to process additional data and implement new algorithms for oil and gas prospecting. Recent work for applications conversion to operate with the IBM 3090-20E/2VF are described and benchmarks results are presented.

## 1. INTRODUÇÃO

As companhias de petróleo, na busca de óleo e gás, vem processando quantidades cada vez maiores de informações. Devido a grande massa de dados a manipular e a complexidade dos algoritmos aplicados a esses dados, surge a necessidade do uso de sistemas de processamento, dotados de alta capacidade computacional.

Para o atendimento às áreas de Exploração, Produção e Perfuração, o Departamento de Exploração (DEPEX) da Petrobrás, vem adquirindo experiência no uso de processadores de alto desempenho desde 1968 com a instalação de um "convolver" (processador de matrizes). Em 1975, foi instalado um "Array Processor" IBM 2938, substituído em 1979 por "Array Processors" IBM 3838 ligados a um computador IBM 308X. Em Setembro de 1987, foi instalado no EDISE/Rio um supercomputador IBM 3090-20E/2VF.

As arquiteturas dos sistemas 3090 com VF e do 3081 com AP-3838 são distintas. A primeira possui um conjunto específico de instruções para manipulação de vetores, que ocorre na própria CPU, sendo que os dados são acessíveis diretamente na memória. Já no 3081 com AP, o cálculo vetorial é executado no "Array Processor", exigindo um tráfego de dados intenso, via canal, entre o processador central e este. A capacidade de processamento máximo teórica de um "Array Processor" IBM 3838 é de 30MFLOPS e o IBM 3090-200 possui uma velocidade escalar um pouco superior a duas vezes a de um IBM 3081.

Por serem arquiteturas distintas, um trabalho de conversão foi desenvolvido e medições foram feitas para avaliação de desempenho.

### 1.1. "Array Processors".

"Array Processors" são processadores periféricos de alto desempenho que aumentam sobremaneira a capacidade de computação do computador ao qual estão acoplados. Normalmente, o uso do AP é feito através de bibliotecas de rotinas codificadas em linguagem de máquina que podem ser chamadas por programas Fortran. A programação para utilização dos recursos de processamento vetorial do AP é bastante complexa pois além da preparação de tarefas especializadas que serão transmitidas para o AP, requer também que o programador faça um mapeamento de memória do AP e mantenha um controle das diversas áreas mapeadas.

Uma grande dificuldade na programação surge na manipulação de vetores longos. Devido a limitação de memória real no AP, o programador deve se preocupar em desenvolver técnicas de fracionamento dos vetores em seções para serem processados.

Outra grande desvantagem dos AP's é que eles requerem transferência de dados de/para o computador hospedeiro, o que normalmente se constitui num gargalo porque esta é uma operação demorada. Os dados são levados, via canal, para o AP onde são processados e os resultados devolvidos, via canal, para o computador central.

### 1.2. IBM 3090-20E com Vector Facility.

A seguinte configuração está instalada no DEPEX:

IBM 3090-20E: 64 Mb de Memória Real, 128 Mb de Memória Expandida e 2 "Vector Facilities".

O IBM 3090-20E tem dois processadores, sendo possível realizar 2 tarefas em paralelo. O 'Vector Facility' (VF) é um dispositivo opcional que pode ser adicionado a cada um dos processadores do 3090 para fornecer alto desempenho no processamento de vetores.

O tamanho do registrador de vetores de um "vector facility" é de 128 elementos. Há 16 registradores de precisão simples ou 8 de precisão dupla. Deste modo o total de memória de registradores vetoriais de cada "vector facility" é 8 Kbytes.

O conjunto de instruções do 3090 foi estendido com 171 novas Instruções Vetoriais executadas no VF. Estas instruções incluem operações matemáticas com Inteiros (Ponto Fixo - 32 bits) e Reais (Ponto Flutuante) de precisão simples (32 bits) e dupla (64 bits) além de operações lógicas em operandos binários. Instruções compostas permitem que duas operações aritméticas sejam feitas por uma mesma instrução vetorial (MULTIPLY AND ADD, MULTIPLY AND SUBTRACT, MULTIPLY AND ACCUMULATE).

Como há dois "pipelines" e cada um pode produzir um resultado por ciclo (tempo de ciclo de 17,2 nanosegundos), teoricamente, o máximo de desempenho do 3090/200/2VF é 232 MFLOPS.

Na prática o uso do VF permite que operações sobre vetores sejam realizadas de 2 a 8 vezes mais rápido que a operação de modo escalar, podendo-se atingir resultados maiores em certos casos.

A memória expandida serve como uma extensão da memória central do 3090, permitindo melhoria de desempenho do sistema. Benefícios potenciais de "THROUGHPUT" podem ser obtidos pela substituição de operações assíncronas de I/O para paginação, por movimentações síncronas, extremamente rápidas, para a memória expandida.

A possibilidade de utilização de memória virtual até 2 gigabytes associada a paginação para memória expandida tem permitido que se tenha reduções significativas no "Elapsed Time" em várias aplicações.

## 2. VETORIZAÇÃO

### 2.1. Software Utilizado.

O suporte de software disponível para os processadores vetoriais é extremamente importante pois é fundamental que os programas possam utilizar todo o potencial existente na máquina. Além disso, ferramentas de análise de desempenho são necessárias para avaliação dos pontos críticos de performance dos programas, de forma a permitir um trabalho otimizado no processo de vetorização.

Software disponível no DEPEX para vetorização e análise de desempenho:

\* MVS/XA.

\* IBM VS Fortran Versão 2 com DEBUG Interativo.

\* IBM ESSL - Engineering and Scientific Subroutine Library.

\* Problem Program Evaluator - da Boole & Babbage.

\* MVS/XA.

O MVS/XA suporta múltiplos espaços de endereçamento de até 2 gigabytes e capacidade de multiprocessamento.

\* IBM VS Fortran V2.

Suporta o nível de linguagem 77, incluindo extensões IBM. Entre as diversas características presentes neste compilador, algumas são de grande ajuda para o desenvolvimento de aplicações vetorizadas:

Níveis de vetorização.

Três níveis de vetorização podem ser selecionados por uma opção de compilação:

\* Vector(level (0))-Não vetorizar.

\* Vector(level (1))-Vetorização "Loop a Loop".

\* Vector(level (2)) - Vetorização por "statement".

Neste último nível, o compilador faz uma separação do loop em dois, um com os "statements" vetorizáveis e o outro não.

Diretivas.

O uso de diretivas fornece ao compilador informações sobre o código, permitindo uma análise mais precisa no processo de vetorização.

Multitasking Facility - MTF.

Pelo uso da Multitasking Facility (MTF) um programa Fortran pode executar código escalar ou vetorial em mais de um processador. A MTF faz parte da biblioteca do VS Fortran V2 e é ativada via "CALL".

Inter-Compilation Analysis - ICA.

Esta opção de compilação solicita ao compilador uma análise extensa das ligações entre unidades do programa. O compilador produz uma tabela de referência cruzada indicando conflitos entre: argumentos, tamanho de blocos COMMON, nomes externos, tipos de funções, etc., permitindo uma agilização do processo de depuração.

DEBUG Interativo - IAD.

Utilizado em tempo de execução, em "batch" ou interativamente, para fazer depuração e análise de desempenho de programas otimizados/vetorizados.

Apresenta os resultados de forma numérica e/ou gráfica, informando o número de vezes e

que cada comando é executado e a porcentagem do tempo de execução gasto em cada subrotina.

Quando executado em terminal, possui recursos como alterações de cor e animação. Esta permite ao programador visualizar a execução dos comandos de seu programa, verificar valores das variáveis, etc.

\* IBM ESSL

A ESSL é uma biblioteca de 233 rotinas matemática, extremamente otimizadas para se obter vantagens da arquitetura do 3090 com VF. Possui aplicações em diversos campos de engenharia e ciência, suportando as seguintes áreas de computação:

- \* Álgebra Linear.
- \* Operações com Matrizes.
- \* Equações Algébricas Lineares.
- \* Eigensystem Analysis.
- \* Processamento de Sinal.
- \* Geração de Números Aleatórios.
- \* Sorting and Searching.
- \* Interpolação.
- \* Quadratura Numérica.

\* P.P.E. da Boole & Babbage.

O Analisador de desempenho da Boole & Babbage permite identificar pontos críticos.

2.2. Metodologia usada para Vetorização de Programas.

No processo de vetorização de programas, algumas etapas foram consideradas de forma a permitir uma maior otimização, tanto do tempo do analista/programador quanto no desempenho das aplicações.

As seguintes etapas têm sido observadas:

a. Selecionar uma ou mais massas de dados representativas dos processos que o programa simula. Um cuidado especial deve ser dado a esta etapa, pois os pontos críticos de desempenho podem alterar em função de parâmetros de entrada.

b. Executar o programa sob o controle de um analisador de desempenho para a identificação de seus pontos críticos. Os programas "Interactive Debug" do VS Fortran V2 da IBM e "Problem Program Evaluator" da Boole and Babbage têm sido usados com o seguinte critério:

IAD - usado quando programas/sub-rotinas são escritos somente em Fortran, pois a forma como os resultados são apresentados permite a identificação dos comandos críticos no programa fonte.

PPE - quando existem programas sub-rotinas em outra linguagem ou quando se deseja obter uma análise mais detalhada dos recursos utilizados pela aplicação (CPU, I/O, "Waits", etc.).

c. Uma vez identificados os pontos críticos, uma análise é feita quanto a possibilidade de vetorização/otimização, métodos de solução e possibilidades de utilização de rotinas da ESSL.

Esta análise permite também fazer uma estimativa dos ganhos que poderão ser obtidos com as modificações identificadas.

d. Introduzir, controladamente, modificações de forma a se ter uma avaliação do ganho obtido e procurar remover inibidores de vetorização, baseado nas informações dos relatórios do compilador.

Este processo deve ser repetido para todas as possibilidades e uma nova análise de desempenho deverá ser feita para identificação de novos pontos críticos.

e. Procurar analisar a possibilidade de uma reestruturação na lógica e/ou na forma de organização dos dados, de modo a permitir a utilização de métodos de solução que melhor utilizem o "Vector Facility".

### 3. RESULTADOS

3.1. "Benchmark" - 3081 com AP 3838 x 3090 com VF.

A diversidade de arquiteturas e formas de paralelismo que existem em todos os supercomputadores torna difícil a medição e comparação destes equipamentos. Para a execução de testes comparativos, geralmente escolhe-se um conjunto representativo de aplicações e mede-se o tempo de CPU, execução, etc. Este processo é denominado "benchmark". A dificuldade reside em selecionar apropriadamente um teste que realmente represente a carga do CPD.

Dado a diversidade de aplicações em uso, foi necessário delimitar o conjunto de programas a ser medido em cada sistema. Pela sua representatividade no processamento global do CPD do DEPEX, foram selecionados os programas de Deconvolução (convolução), Migração (interpolação e soma de vetores), Empilhamento (soma de vetores) e o de Análise de Velocidades (interpolação) e que representam respectivamente 54% e 78% do processamento sísmico aplicado às linhas terrestres e marítimas (que representam 80% da carga total do CPD).

As medidas foram feitas com a execução de uma cópia de cada programa (execução isolada) e também com a execução de várias cópias (execução conjunta) [1]. Esta última medida pretende determinar a quantidade máxima possível de trabalho em cada sistema para cada tipo de serviço. Os resultados obtidos estão mostrados na tabela a seguir:

Tabela 1. Resultados do benchmark 3081 com 2 x 3838 x 3090 com 2VF

Programas	EXECUÇÃO			
	ISOLADA		CONJUNTA	
	traços/ seg.	3090/ 3081	traços/ seg.	3090/ 3081
Decon				
3090	17,36	2,7	84,91	3,7
3081	6,50		23,28	
Migra				
3090	0,56	9,3	1,17	6,9
3081	0,06		0,17	
Empilha- mento:				
3090	236,38	7,9	716,27	8,0
3081	29,84		89,63	
Análise de Velo- cidade:				
3090	37,99	4,7	89,24	4,4
3081	8,10		20,24	

Traços processados por segundo em cada sistema e relação entre os dois sistemas.

Ponderando-se o ganho obtido em cada programa com a porcentagem de sua participação no processamento global, chegou-se a um ganho médio de 6,8 vezes no processamento terrestre e 6,3 vezes para o marítimo.

### 3.2. Vetorização de Funções Sísmicas.

Foram feitos testes em Funções Sísmicas Modulares da PETROBRÁS [2], com a finalidade de determinar o grau de vetorização e o ganho em tempo de CPU obtido pela vetorização dessas rotinas.

As seguintes Funções Modulares foram avaliadas, fazendo uso das versões escalar e vetorial da biblioteca ESSL.

#### QFCOE - Filtro de Coerência.

Métodos matemáticos: Interpolação, Convolução, Soma de Vetores.

#### QFVEL - Filtro de Velocidade.

Métodos matemáticos: Convolução, Soma de Vetores.

#### QDARC - Determinação de Anomalias em Registro de Campo.

Métodos matemáticos: Soma de Vetores, Convolução, Funções Estatísticas.

#### QIMPE - Impedância Acústica.

Métodos matemáticos: FFT Forward Complex, FFT Inverse Real.

A tabela a seguir mostra os resultados obtidos para as versões escalares e vetorizadas dessas funções.

Tabela 2. Vetorização de Funções Sísmicas

Função	Versão ESSL Escalar	Versão ESSL Vetorial		Relação ESC/VET
	Tempo CPU/seg.	Total	Vetor.	
QFCOE	1056	131	113	8,1
QFVEL	1868	214	174	8,7
QDARC	240	32	25	7,5
QIMPE	336	47	31	7,1

### 3.3. Outros testes.

Visando o treinamento do pessoal, além de seminários e cursos específicos, foi realizado um trabalho de vetorização com quatro aplicações de diferentes áreas [3]. Foi utilizada a metodologia descrita no item 2.

Na aplicação de "Estimativa de Malhas" a compilação vetorial indiscriminada (intencional) de todo o programa aumentou o tempo de CPU em 4,7%. Após a vetorização seletiva, a modificação de algoritmo e outras pequenas alterações conseguiu-se um ganho de 1,83 (de 73,84 seg. para 40,28 seg.).

Em um programa de simulação pelo método de Monte Carlo, a seleção randômica de elementos de um conjunto tridimensional de 2.880.000 elementos, não foi favorável ao processo de vetorização. Após todo o esforço de análise, o máximo obtido foi um ganho de apenas 1,6%.

Na aplicação "Análise Não Linear de Estruturas" verificou-se que a melhor solução para este programa será um novo projeto de algoritmo, com total reestruturação dos dados.

Para um modelo de simulação de engenharia de reservatório de hidrocarbonetos com 532 blocos (solução de sistemas de equações lineares) conseguiu-se com a vetorização a redução de 75% (quatro vezes) no tempo de CPU e 87% (7,4 vezes) no tempo de execução.

## 4. CONCLUSÕES

Importantes observações têm sido verificadas na experiência da PETROBRÁS.

a) O uso de uma metodologia de vetorização e avaliação é importante para atingir resultados em prazos aceitáveis;

b) Qualquer tentativa para obtenção de alto desempenho de programas está ainda intimamente ligada à arquitetura da máquina;

c) Nem todas as aplicações são apropriadas para atingir um bom desempenho com vetorização;

d) A vetorização indiscriminada de todas as rotinas, através do compilador Fortran, sem levar em conta a análise dos pontos críticos, não produz os melhores resultados, podendo em determinados casos piorar o desempenho da aplicação;

e) O uso da biblioteca ESSL é fundamental para se conseguir altos níveis de "speed-up ratio" no "vector facility" do 3090;

f) A vetorização de programas propicia uma oportunidade de se rever antigos conceitos de otimização de programas.

## 5. AGRADECIMENTOS

Ao Tamanini, Paulo F. Santos, Ana Lúcia Gomes e Ivan Pedroza pela revisão e sugestões. À Dna. Jeanice Gedeon pela montagem final.

## REFERÊNCIAS

- [1] Informe Geofísico n.76 - Medidas de desempenho dos sistemas: IBM 3090-VF e 3081-K - Array Processor - Janeiro de 88 PETROBRÁS/DEPEX.
- [2] Relatório de Vetorização - Março de 1988 DITREX/SESCEF.
- [3] Relatório de Vetorização de Programas de Usuários - PETROBRÁS/DEPEX/DISEP/SEAP.