# 平成27年度オペレーティングシステム期末試験

# 2016年1月19日(火)

- 問題は3問
- この冊子は, 表紙 1 ページ (このページ), 問題 2-8 ページからなる
- 解答用紙は1枚. おもてとうらの両面あるので注意すること.
- 各問題の解答は所定の解答欄に書くこと.

### 1

巨大なメモリを搭載するコンピュータ上で、以下の AO~A4 までの 5 種類プログラムを走らせることを考える.

#### A0:

```
int main() {
   char * a = malloc(A);
   /* --- */
}
```

#### **A1:**

```
int main() {
  int fd = open(file, O_RDONLY);
  char * a = malloc(A);
  ssize_t r = read(fd, a, A);
  read_array(a, B); /* 中身は後述 */
  /* --- */
}
```

#### **A2**:

```
int main() {
    int fd = open(file, O_RDONLY);
    char * a = mmap(0, A, PROT_READ|PROT_WRITE, MAP_PRIVATE, fd, 0);
    read_array(a, B); /* 中身は後述 */
    /* --- */
}
```

#### **A3**:

```
int main() {
    int fd = open(file, O_RDONLY);
    char * a = mmap(0, A, PROT_READ|PROT_WRITE, MAP_PRIVATE, fd, 0);
    write_array(a, B); /* 中身は後述 */
    /* --- */
}
```

#### **A4**:

```
int main() {
    int fd = open(file, O_RDWR);
    char * a = mmap(0, A, PROT_READ|PROT_WRITE, MAP_SHARED, fd, 0);
    write_array(a, B); /* 中身は後述 */
    /* --- */
}
```

ここで read\_array(a, m) および write\_array(a, m) は、それぞれ a から始まる m バイトを読み出すまたは書き込む関数で、以下に相当する動作をする.

```
void read_array(char * a, ssize_t m) {
    for (ssize_t i = 0; i < m; i++) {
        a[i];
    }
}

void write_array(char * a, ssize_t m) {
    for (ssize_t i = 0; i < m; i++) {
        a[i] = 'x';
    }
}</pre>
```

(1) AO を走らせた時消費する物理メモリの量 (おおよその値) を答えよ.答えはプログラム中のパラメータ A, B を用いた式で表せ.ただし A や B はある程度 (少なくとも数十 M) の大きな値である.A1~A4 についても同様に答えよ.

以下を仮定してよい.

- システム関数呼び出し (open, malloc, mmap, read) は全て成功する.
- read は要求したバイト数分のデータを実際に読み出す.
- 関数呼び出しやメモリアクセスが実際の計算には不要だからといって、コンパイラの最適化で除去されることはない (プログラムはあくまで書かれた通りのことを実行する).
- 上記プログラム中で明示的に行われるメモリ割当量に比べて、それ以外に必要な (例えばプログラムのコードやスレッドのスタックのための) メモリ量は小さいと考え、無視する (0 とする).

なお,消費されるメモリ量は,走らせているプログラムが/\* --- \*/の行に達した時点での消費量とする.

- (2) A0 を P 個同時に立ち上げるとどうか? 答えは A, B, P を用いた式で表せ. 消費されるメモリ量は,P 個全でのプロセスが/\* --- \*/に達し (そこでしばらく停止し) たところでの消費量とする. A1~A4 についても同様に答えよ.
- (3) 以下は OS がメモリ管理のために用いてる技法であるが、それぞれどのようなもので、どのような時にどのような利点があるか、簡潔に説明せよ.
  - (a) 要求時ページング
  - (b) Copy on write
- (4) (2) の A2 および A3 の物理メモリ使用量の結果を、要求時ページング、copy on write という言葉を、関連があれば適宜使いながら、それぞれ説明せよ.

# 2

以下の関数  $f0\sim f3$  はいずれも,あるメモリ領域に N 回 1 を足し,終了後にその値を表示するという関数であり,対象となるメモリ領域だけが異なる (for 文以降はほとんど共通: 重要部に下線が引いてある).

f0:

```
int x = 0;
void * f0(void * _) {
   for (int i = 0; i < N; i++) {
      x = x + 1;
   }
   printf("it's %d\n", x);
   return 0;
}</pre>
```

f1:

```
void * f1(void * _) {
    int * p = (int *)calloc(1, sizeof(int));
    for (int i = 0; i < N; i++) {
        *p = *p + 1;
    }
    printf("it's %d\n", *p);
    return 0;
}</pre>
```

f2:

```
void * f2(void * arg) {
    int * p = ((arg_t *)arg)->m;
    for (int i = 0; i < N; i++) {
        *p = *p + 1;
    }
    printf("it's %d\n", *p);
    return 0;
}</pre>
```

f3:

```
void * f3(void * _) {
   int fd = open("/tmp/zero_file", O_RDWR);
   int * p = mmap(0, sizeof(int), PROT_READ|PROT_WRITE, MAP_SHARED, fd, 0);
   for (int i = 0; i < N; i++) {
      *p = *p + 1;
   }
   printf("it's %d\n", *p);
   return 0;
}</pre>
```

ただし、f2における arg\_t は以下のような構造体である.

```
typedef struct {
  int * m;
} arg_t;
```

また, f3 における/tmp/zero\_file は事前に作られ, 中身はすべて 0 で埋まった 4 (sizeof(int)) バイトのファイルであるとする.

 $f0\sim f3$  それぞれの関数を,P 個のスレッドもしくは P 個のプロセスで並行に走らせる.例えば f0 を P 個のスレッドで走らせる場合,

```
void run_threads() {
     arg_t args[P];
     int * m = (int *)calloc(1, sizeof(int));
     for (int i = 0; i < P; i++) {
       args[i].m = m;
     }
     pthread_t tids[P];
     for (int i = 0; i < P; i++) {
       pthread_create(&tids[i], 0, f0, &args[i]);
     }
     for (int i = 0; i < P; i++) {
11
       pthread_join(tids[i], 0);
12
     }
   }
14
```

のようにする. f0 を P 個のプロセスで走らせる場合,

```
void run_procs() {
     arg_t args[P];
     int * m = (int *)calloc(1, sizeof(int));
     for (int i = 0; i < P; i++) {
       args[i].m = m;
     }
     pid_t pids[P];
     for (int i = 0; i < P; i++) {
       pids[i] = fork();
       if (pids[i] == 0) {
10
         f0(&args[i]);
11
         exit(0);
       }
13
     }
14
     for (int i = 0; i < P; i++) {
       waitpid(pids[i], 0, 0);
16
     }
```

のようにする (2-6 行目は run\_threads と run\_procs で共通である).

f1~f3 を走らせる場合は、それぞれの中の f0 と書かれた部分 (run\_threads の 9 行目および run\_procs の 11 行目) を、適切な関数名に置き換えるだけである。これで都合 2 (run\_threads, run\_procs) × 4 (f0~f3) = 8 通りのプログラムができたことになる。各スレッド (ないしプロセス) が、"it's x" という表示をするので、いずれのプログラムもそのような行を P 個表示することになる。

以下ではN = 10000000, P = 2とする.以下の問いに答えよ.

(1) 8 通りのプログラムそれぞれについて、その出力としてあり得るものを、以下の中から全て選び、記号で答えよ.

(ア)

```
it's 8193611
it's 10000000
```

(1)

```
it's 8193611
it's 10242445
```

(ウ)

```
it's 10000000
it's 10000000
```

(工)

```
it's 10000000
it's 10242445
```

(オ)

```
it's 10193611
it's 10242445
```

(カ)

```
it's 19706152
it's 20000000
```

(キ)

```
it's 19706152
it's 20000001
```

(ク)

```
it's 20000000
it's 20000000
```

解答は、解答用紙の表で、あり得るケースにチェック(√)を入れて答えること.

(例)

		(ア)	(イ)	(ウ)	(工)	(オ)	(カ)	(キ)	(ク)
run_threads	fO	✓	✓	✓	✓				
run_threads	f1		✓	✓		✓			✓
:	:								

- (2) 例えば(オ)のような出力が生じうるケースについて、それがどのようにして生じるのかを説明せよ.
- (3) pthread\_mutex\_t 型の大域変数 M を用意し、各繰り返しの中で1を足す操作の前後を pthread\_mutex\_lock(&M) と、pthread\_mutex\_unlock(&M) ではさむことにする.

具体的には、f0であれば以下のように変更する. 下線部 が変更した部分.

```
pthread_mutex_t M;
int x = 0;

void * f0(void * _) {
   for (int i = 0; i < N; i++) {
        pthread_mutex_lock(&M);
        x = x + 1;
        pthread_mutex_unlock(&M);
}

printf("it's %d\n", x);
return 0;
}</pre>
```

 $f1\sim f3$  についてもほぼ同じ変更を施す (唯一の違いは、pthread\_mutex\_lock(&M), pthread\_mutex\_unlock(&M) の呼び出しではさむのが、\*p = \*p + 1; である点).

このようにしてできたプログラム達の出力はどうなるか? (1) と同じ形式で答えよ.

- (4) 多くのプロセッサに備わる compare-and-swap 命令について, その動作を説明せよ.
- (5) compare-and-swap 命令を用いて、pthread\_mutex\_lock、pthread\_mutex\_unlock で達成したのと同じよう なことを達成することができる. 具体的にどうすればよいか述べよ. 答え方としては、f0 の x = x + 1; ま たは  $f1 \sim f3$  の\*p = \*p + 1; の行を、どのように変更すればよいかを記せ (どちらを使っても良い).

少年は指導教員から、「自分のノート PC のハードディスクの性能を測れ」と言われた。そこで少年は、自分の PC (これはやや旧式で、SSD ではなくハードディスクを搭載している) 上で 1GB のファイルを作り、それを読みだすのにかかる時間を測定した。

少年はこの結果を持って、「1GB が 0.149 秒で読めたとあるから、ハードディスクの読み出し性能は 1GB/0.149 秒  $\approx 6.7$ GB/秒くらいです、ついでに言うと、書き込み性能は dd の結果から、354MB/秒くらいのようですね」と報告したところ、指導教員に  $\underline{(1)}$  アホかと言われ てしまった。そして、「ハードディスクの性能が測りたいなら、あることに気をつけて測定しないとダメだ」と言われた。そこで (2) きちんと気をつけて測定 したところ、

```
$ time cat large > /dev/null
real 0m7.902s
user 0m0.007s
sys 0m0.212s
```

という結果が得られ、約 120MB/秒程度の性能となり、指導教員にはその数字ならあり得ると言われた.

次に、「1GB のファイルが 7.902 秒で読み出せることはわかった。 じゃあ,1GB のファイル中の,指定された 1KB の部分 (だけ) を読めと言われたら,それにかかる時間はどのくらいか」と聞かれた. 少年は調べてきますと言いつつ、「1GB が 7.902 秒だから 1KB にかかる時間は,実験などしなくても簡単に計算できる」と思い,

7.902 秒 × 1KB/1GB  $\approx 7.902 \times 10^{-6} \approx 8 \mu$  秒

さも実験をしたようなふりをしつつ「約8マイクロ秒でした」と答えたところ,「ウソつけ」と言われしまった. (3)「カタログに載っているこの数字を見れば,数ミリ秒より速いはずがないんだよ」とのことであった.

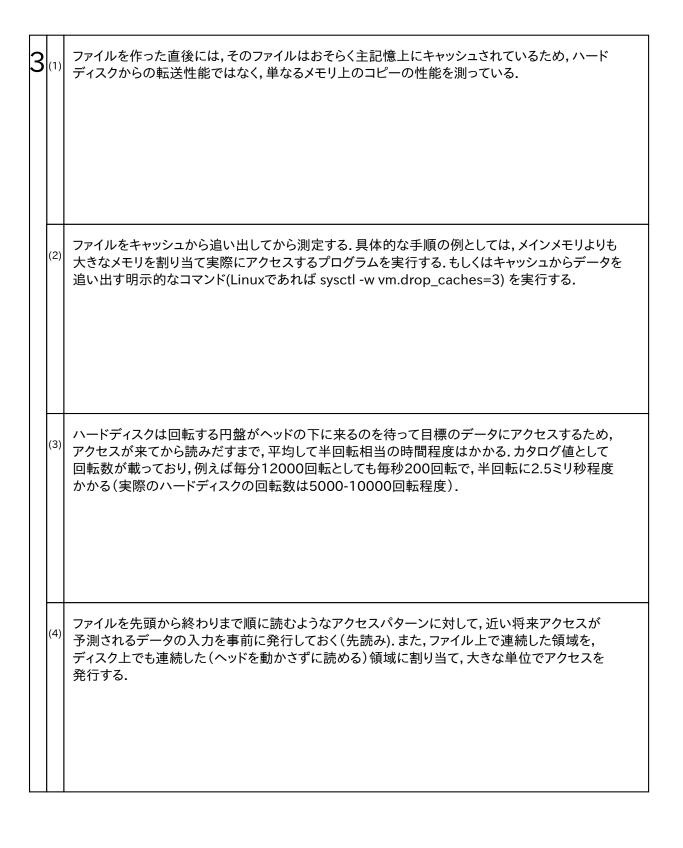
少年は、OSが行っているファイル読み書き性能が、いくつもの (4) ソフトウェア的な工夫によって達成されていることを学んだのであった。

以下の問いに答えよ.

- (1) 下線部 (1) について, 少年の測定はなぜ間違いなのか, 説明せよ. 6.7GB/秒という数字がハードディスクの性能でないとしたら, いったいなぜこんな数字が出たのか?
- (2) 下線部(2)について,正しく測定するための手順(の例)を述べよ.
- (3) 下線部(3) について、カタログに載っているどのような数値を元に、「数ミリ秒より速いはずがない」と判断できるのか?
- (4) 下線部 (4) について, 仮に 1KB 読むのに数ミリ秒はかかるのだとしたら (仮に 1 ミリ秒かかるとする) これを 比例させると, 1GB 読むのには 1000 秒はかかる計算になる. それが実際には 7.9 秒で読み出せているのはど のような工夫を OS がしているからなのか?

問題は以上である

所属学科						氏名											
				+-				1									
	(1)			(2				• \									
1		A1 2A	(A △)		_			1) (O(A	P)∆)								
ı		A2 B	(2B △)			2 B		(2B △)									
		A3 2B (	B, A+B △)					1) (O(B	P) ∆)								
		A4 B	(2B △)			4 B		(2B ∆)									
	(a) sbrkやmmapなどのメモリ割り当てAPI呼び出し時には物理メモリを割り当てず、その後実際に																
	│									f点でのB	侍間						
	(3)	を短縮で															
	<u> </u>	(b) 複数プロ	セスがfoi	kや	PRI	/ATE	なmm	iapのí	洁果同	一ペ-	ージの	コピー	·を持 <sup>·</sup>	つこと	になっ	った場合	,
		実際に書	き込みが	行扎	っれる	れるまで物理ページを共有し、書き込みが起きた時点でコピーを作る.									を作る	それら	
		のページ	が実際に	実際には余り変更されない場合,forkやmmapの高速化,物理メモリの節約になる										,			
	A2 要求時ページングの効果によりファイルの大きさによらず実際にreadした大きさ(配列Aので 込まれた分)だけしか物理メモリは割り当てられない、またmmapした領域に書き込まないの																
										よいので	,						
	(4)		,									•					
	(-)	A3 要求時ペ															きき
		込まれた															
		各プロセ	,									175 111	Парс	, C 1-50	>,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,		.,
	_		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	,	,,			1		(1)		1 (1)	1				
	(1)			Ш	(ア)	(イ)	(ウ)	(工)	(オ)	(力)	(+)	(ク)					
		run_threa		<	7	<b>✓</b>		1	1	<b>✓</b>	lacksquare	/	>				
2		run_threa		Ш			//										
		run_threa			1	1	1	1	1	✓		~	5				
		run_threa			<b>~</b>	✓	1	1	1	/		_					
		run_procs					1										
		run_procs		Ш			1										
		run_procs		Ш			•										
		run_procs	f3	<	1	1	<b>/</b> /	✓	1	✓	ackslash J	/	>				
	(2)	一つのスレ	ッドT0が	xを	読んご	でから	書くま	での間	引に. 州	<sup>1</sup> のス	レッドフ	1がx	に書し	いてい	ると	T1による	)
	(2)	インクリメン												-	,		
									,				1				
	(3)		1 60		(ア)	(1)	(ウ)	( <u>T</u> )	(才)	(力)	(+)	(ク)					
		Tun_tinea			$/\!$	<u> </u>	//	<u> </u>	$\rightarrow$		$\!$	/	2				
		run_threa		Щ	<sup>7</sup>	$\longrightarrow$	//	<u> </u>	$\longrightarrow$								
		run_threa		$\bot \! \! \! \! \! \! \! \! \! \! \! \! \! \! \! \! \! \! \!$		\				1		1					
		run_threa		Ш						~		1					
		run_procs		Ц			1										
		run_procs		Ц			1										
		run_procs		Ш		/_	1	<u>\</u>	/								
		run_procs	f3		7	V	1	V	<b>1</b>	1		1	>				
		compare-	and-swa	n n	V V	/ (//.±	レジス	タ)で	C₩Ħ	で聿に	ナばい	かにお	3 4 古	スァン	を不	可分に行	i à
(4) compare-and-swap p, x, y (yはレジスタ)で, C言語で書けばいかに相当することを7 if (*p == x) swap(*p, y)									۰۰۲۲	-1 /1 (C.1.1	٦.						
		Π(β λ	.) 3wap(	J, y,	,												
	(_`	一回のイン	′クリメント	を以	以下の	りように	に行う										
	(5)	<sup>5)</sup> while (1) {															
		t = x; y	= t + 1;														
		compare	e-and-sw	ар	p,t,y	;											
		if $(y == t)$	) break;		-												
		] }	,														



## 解説

1

 $\mathbf{2}$ 

あり得る表示は,

$$\begin{array}{ccccc}
N & \leq & m & \leq & 2N, \\
0 & < & M-m & < & N-1
\end{array} \tag{1}$$

を満たす (m, M) の組全て及びそれだけである.

まず上記を満たす任意の(m, M)の組に対して、実際にそのような値が表示されうることを示す。

 $0 \le a \le N-1$  を満たす任意の a に対して、以下のような実行の履歴があり得る.

_	_		
	$T_0$	$T_1$	x の値
		$R_0$	0
	$R_0$		0
	:		:
	$W_{a-1}$		a
		$W_0$	1
	$R_a$		1
	:		:
†	$W_{N-1}$		N - a + 1
		$R_1$	N-a+1
		:	:
*		$W_{N-1}$	2N-a

よって特に、a=2N-M とすれば、 $T_1$  が M を表示することになる。またこのとき  $T_0$  が最後の書き込みを終了した時点 († の行の終了後) での  $\mathbf x$  の値は -N+M+1 で、 $T_0$  はここから \* までの任意の時点で printf で表示する値を読みうるので、

$$-N+M+1 \le m \le M$$
,

つまり

$$0 \le M - m \le N - 1$$

を満たす任意のmに対して、そのmが $T_0$ によって表示されうることになる。

逆にどのような実行系列で表示される値も必ず式 (1) を満たすことを示そう。実行している二つのスレッドを  $T_0, T_1$  としよう。まず両者とも以下のようなアクション R と W を交互に N 回ずつ,繰り返し実行していることを確認しておく。

$$R: t = x;$$
  
 $W: t = t + 1; x = t; i = i + 1;$ 

ひとつのアクションは共有変数 (x) へのアクセスと一度しか含んでいないので,不可分に実行されるとみなして良い.実際に起こりうる実行系列は,スレッド  $T_0$  が実行する  $R,W,R,W,\ldots$  と  $T_1$  が実行する  $R,W,R,W,\ldots$  を,任意にインターリーブさせたものである.  $^1$ 

キーとなる不変条件は以下である. 各繰り返しの先頭 (R実行直前および W 実行直後) 以下が成り立つ.

$$\min(i_0, i_1) \le \mathbf{x} \le i_0 + i_1 \tag{2}$$

ただし $i_0(i_1)$ は、 $T_0(T_1)$ が持つ変数iの値である.

この条件が、あるスレッドが R;W を一回実行した時に保存されることを示そう.一般性を失うこと無く、 $T_0$  が R,W を一度ずつ実行する系列に注目する.

$$R:$$
  $\mathbf{t} = \mathbf{x};$   $\vdots$   $\vdots$   $T_1$ の任意の命令列  $\vdots$   $\vdots$   $W:$   $\mathbf{t} = \mathbf{t} + \mathbf{1}; \mathbf{x} = \mathbf{t}; \mathbf{i} = \mathbf{i} + \mathbf{1};$ 

という一連のアクションの前後で,式 (2) が保存されることを示す。 $T_1$  の i の値の,R 実行前時点での値を  $i_1$ , W 実行後時点での値を  $i_1'$  と書いて区別する。同様に x の値の,R 実行前時点での値を x, W 実行後時点での値を x' と書いて区別する。

$$\min(i_0, i_1) \le x \le i_0 + i_1 \Rightarrow \min(i_0 + 1, i_1') \le x + 1 \le i_0 + 1 + i_1'$$

RとWの間に,  $T_1$ がWを $n(\geq 0)$ 回実行したとする. すなわち $i_1'=i+n$ . 一般に,

$$\min(x + a, y + b) \le \min(x, y) + \min(a, b)$$

に注意しておく.

左側の不等式は,

$$\min(i_0 + 1, i'_1) = \min(i_0 + 1, i_1 + n) \le \min(i_0, i_1) + \min(1, n) \le \min(i_0, i_1) + 1 \le x + 1$$

だからよい. 右側の不等式は,

$$x + 1 \le i_0 + i_1 + 1 \le i_0 + (i_1 + n) + 1$$

だからよい.