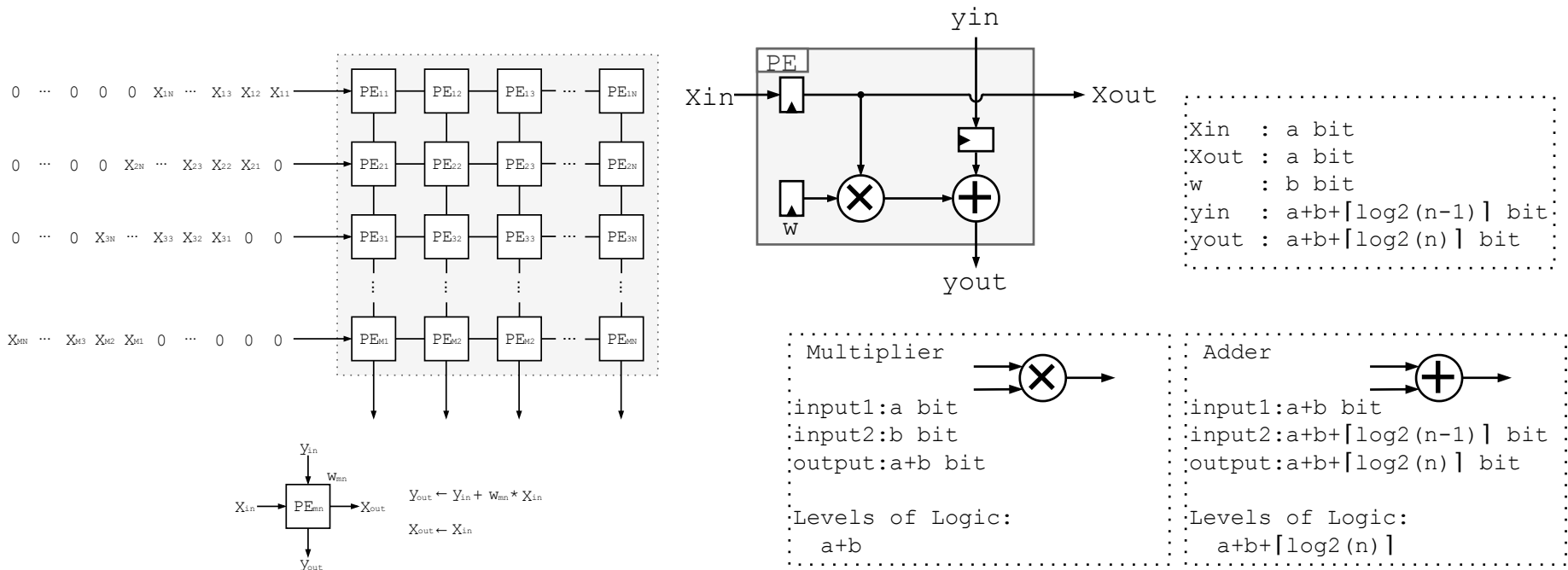


Pipeline Systolic Reconfigurable Fixed-Point Accelerator

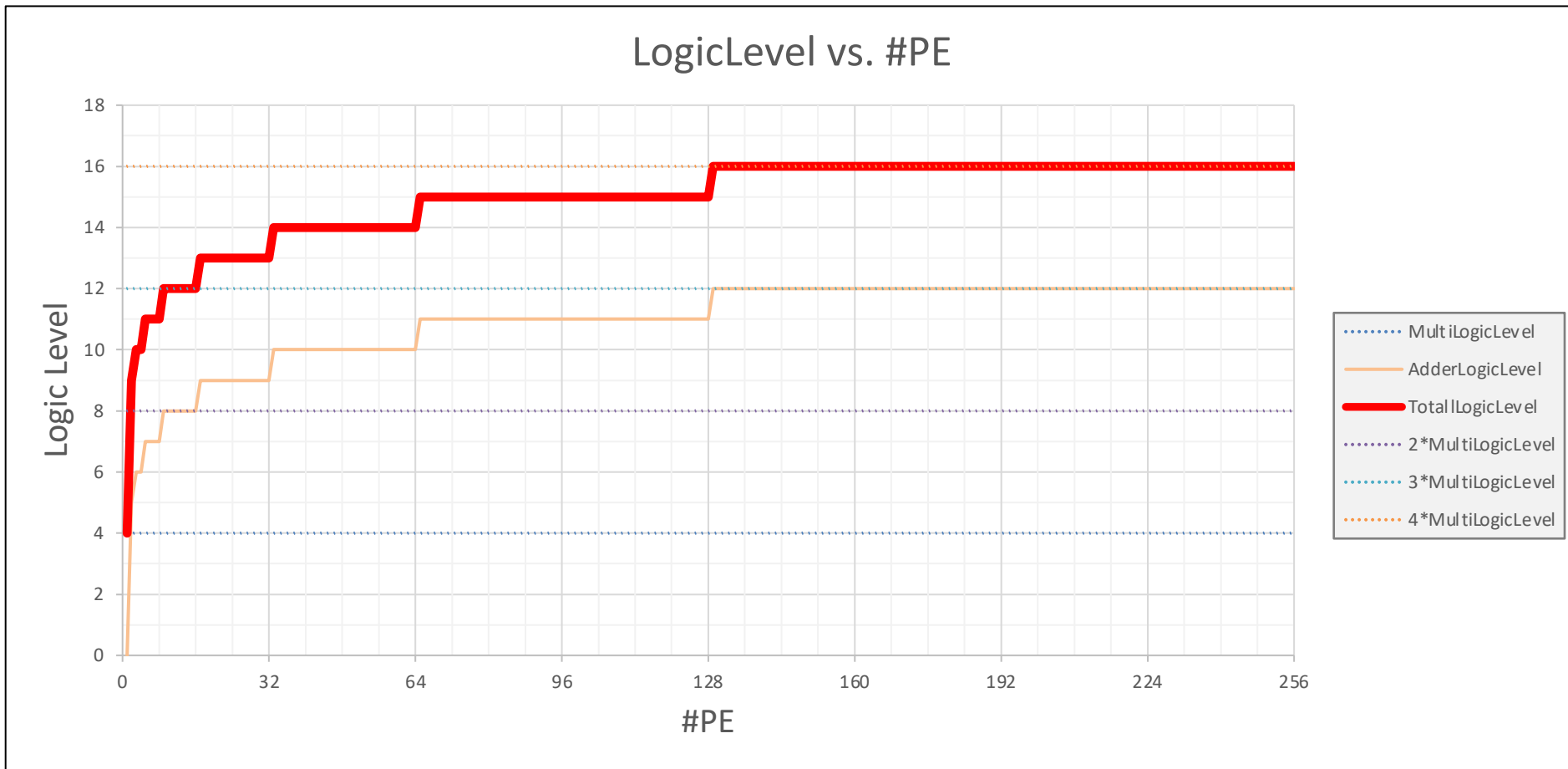
文章计划

Background



- 如上图，脉动阵列的Psum长度会随着多次累加的长度而上升（对数上升），其使得加法器的一个input及output长度的增加，从而增加了Logic-Level，增加延时。
- 混合精度加速器的基础W/X带宽为2bit，而 $\log_2(256)=8$ 。对于低精度脉动阵列而言，这部分增加的延时对乘加速度有很大影响。

Background



- 以2bit*2bit为基本MAC单元的阵列，当阵列宽度达到N=256时，加法的Logic-Level（延时）是乘法的3倍。
- 若对MAC采取分段式运行4-states Pipeline，则能提高4倍的运行速度。

Architecture of PE

```

n >= 1 (从1开始计数)
t = (a+b+[log2(n)]) / (a+b+1)
PEmn has [t] states pipeline
Specially, when [t]=t
PEmn has Cout
    
```

```

Xin : a bit
Xout : a bit
w : b bit
yin : a+b+[log2(n-1)] bit
yout : a+b+[log2(n)] bit

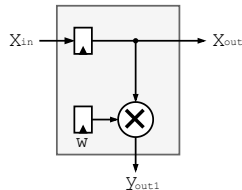
yin1 ... yin[t-1] : a+b bit
yout1 ... yout[t-1] : a+b bit

yin[t] : (a+b+[log2(n-1)])%(a+b) bit
yout[t] : (a+b+[log2(n)])%(a+b) bit

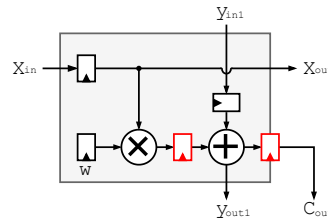
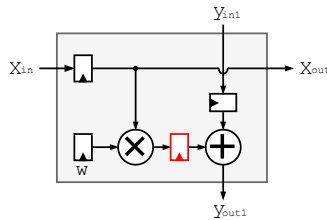
Cout : 1 bit

yin = [ yin[t] ... yin2 yin1 ]
yout = [ Cout yout[t] ... yout2 yout1 ]
    
```

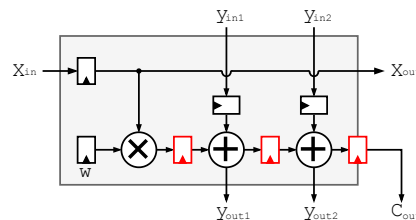
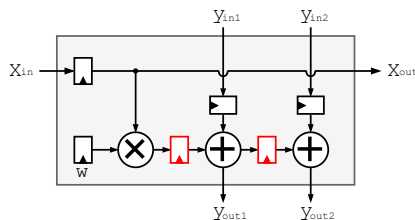
1 states:



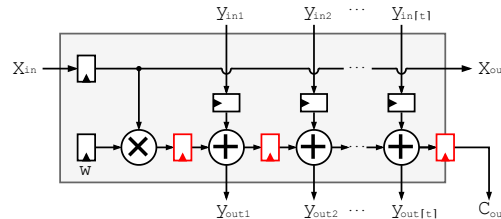
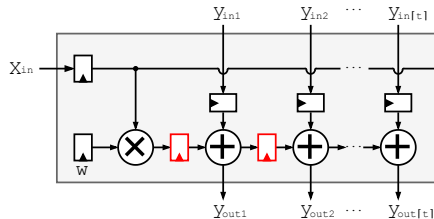
2 states:



3 states:

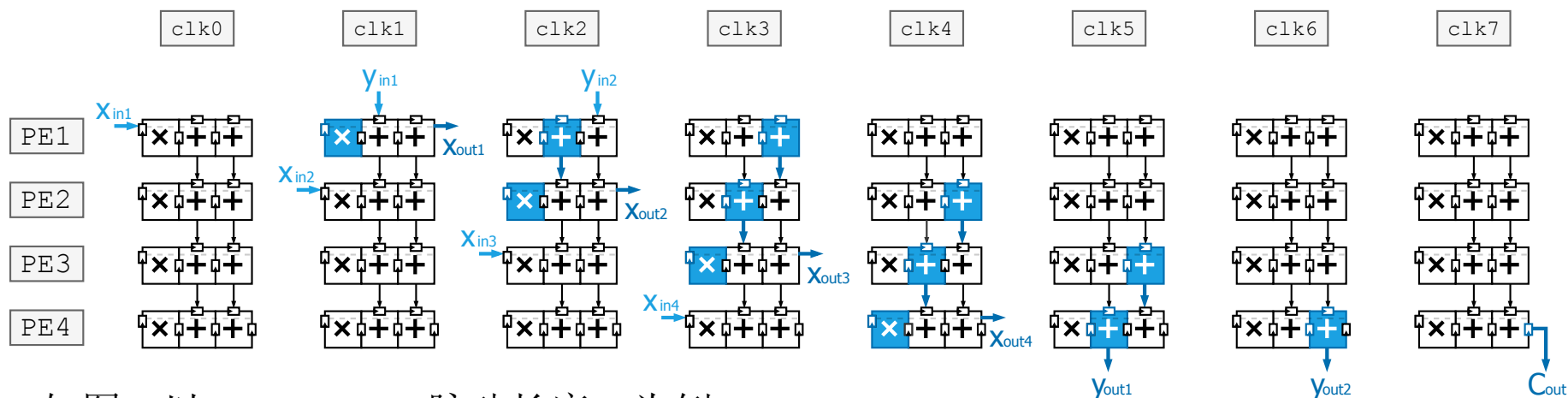


[t+1] states:



- 左图是多级pipeline的设计
- PE内部的pipeline长度依据 Y_{in} 的长度积累
- 阵列中第n排PE的架构相同，每当 $\log_2(n)$ 为整数时，psum位宽+1。
- 位宽没增加 $a+b$ bit，新增 1state 的pipeline
- 乘法器-加法器的寄存器为 $a+b$ bit
- 加法器-加法器的寄存器为 1bit
- 从 2states 开始，每增加一级pipeline，其边际成本 = 1bit 寄存器

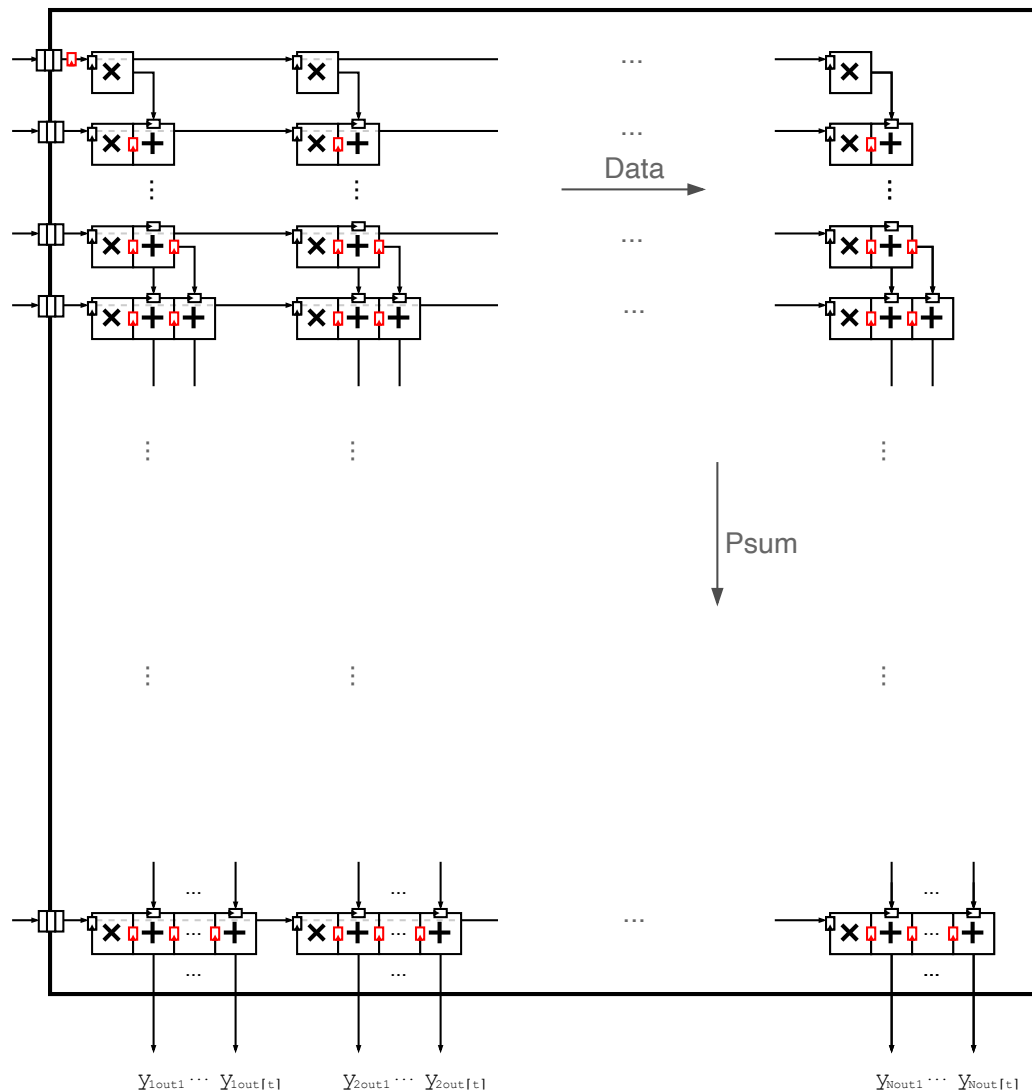
Pipeline + Systolic



如图，以3-states PE，脉动长度=4为例

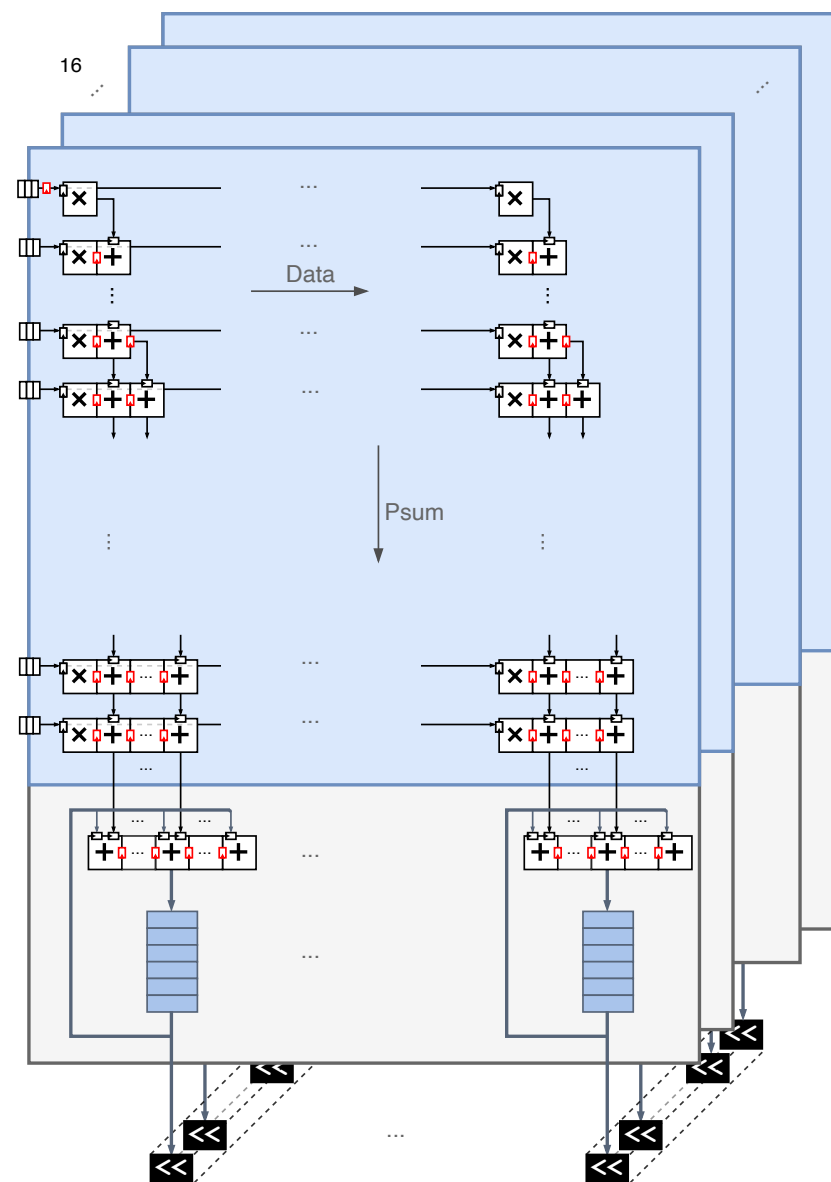
- $X_{in1} - X_{in4}$ 分别在 $clk0 - clk3$ 输入给 PE1 - PE4。
- Y_{in1}, Y_{in2} 分别在 $clk1, clk2$ 输入给 PE1。
- X_{in} 一个 clk 后与 PE 内部的权重 W 相乘后生成的乘积 $X*W$
- 乘积与 Y_{in1} 相加的 C_{out1} 与 Y_{in2} 相加
- 在某些情况（ n 为 2 的指数）下， C_{out1} 与 Y_{in2} 相加会产生 C_{out2}
- Y_1, Y_2 在 PE 中流动，相隔 $1clk$ ，分别
- C_{out2} 在 $clk7$ 时从 PE4 输出

PE Array



- 增设了pipeline的脉动阵列的输入&控制行为与传统脉动阵列相同。
- 由于增设了pipeline，输出数据为多段数据，相邻数据段相隔1clk延迟输出。

Architecture of Reconfigurable Accelerator



- 加速器采用TPU架构，采用先脉动再移位的思想。在加速器层次进行多精度可重构设计。
- 每层阵列的精度为2bit*2bit，共16层阵列。可支持2-8bit混合精度数据进行矩阵乘法计算。
- 阵列模块后接累加模块&移位模块。累加也采用pipeline加法。

Evaluation and Comparison

PE Comparison	(1) Bit fusion (2) Bit serial
Accelerator Comparison	(1) TPU (2) Eyeriss
Evaluation	Vgg / Alexnet
Evaluation	(1) Power (2) Area (3) Latency/speed (4) Performance (Gops)/Efficiency (Gops/W) (Mixed precision)

投稿相关信息

时间	周俊卓	满昌海	罗君益
3.29-4.11	PE单元实现&测试		
4.12-4.25	PE阵列实现&测试		
4.26-5.16	多层PE阵列多精度实现&测试		SoC系统搭建
5.17-5.23	加速器总线接口设计	测试模型训练&转化	
5.24-6.6	加速器SoC整体测试（仿真）		文章理论部分
6.7-6.20	FPGA综合+测试		
6.21-7.26	完善文章理论部分 编写实验部分		

	abstract deadline	deadline	length
ASPDAC	2021.7.26	2021.7.26	8 pages
DATE	2021.9.14	2021.9.21	6 pages