

Mapping Adaptive Particle Filters to Heterogeneous Reconfigurable Systems

Thomas C.P. Chau, Department of Computing, Imperial College London, UK

Xinyu Niu, Department of Computing, Imperial College London, UK

Alison Eele, Department of Engineering, University of Cambridge, UK

Jan Maciejowski, Department of Engineering, University of Cambridge, UK

Peter Y.K. Cheung, Department of Electrical and Electronic Engineering, Imperial College London, UK

Wayne Luk, Department of Computing, Imperial College London, UK

This paper presents an approach for mapping real-time applications based on particle filters to heterogeneous reconfigurable systems, which typically consist of multiple FPGAs and CPUs. A method is proposed to adapt the number of particles dynamically and to utilise run-time reconfigurability of FPGAs for reduced power and energy consumption. A data compression scheme is employed to reduce communication overhead between FPGAs and CPUs. A mobile robot localisation and tracking application is developed to illustrate our approach. Experimental results show that the proposed adaptive particle filter can reduce up to 99% of computation time. Using run-time reconfiguration, we achieve 25-34% reduction in idle power. A 1U system with four FPGAs is up to 169 times faster than a single-core CPU and 41 times faster than a 1U CPU server with 12 cores. It is also estimated to be 3 times faster than a system with 4 GPUs.

Categories and Subject Descriptors: C.1.3 [Processor Architectures]: Other Architecture Styles—*heterogeneous (hybrid) systems*; C.3 [Special-purpose and Application-based Systems]: Real-time and embedded systems

General Terms: Algorithms, Design, Performance

Additional Key Words and Phrases: Particle filters, Sequential Monte Carlo, Reconfigurable systems, FPGAs, Run-time reconfiguration

ACM Reference Format:

Thomas C.P. Chau, Xinyu Niu, Alison Eele, Jan Maciejowski, Peter Y.K. Cheung and Wayne Luk, 2014. Mapping Adaptive Particle Filters to Heterogeneous Reconfigurable Systems *ACM Trans. Reconfig. Tech. Syst.* 1, 1, Article 1 (March 2014), 17 pages.
DOI: <http://dx.doi.org/10.1145/0000000.0000000>

1. INTRODUCTION

Particle filter (PF), also known as sequential Monte Carlo (SMC) method, is a statistical technique for dynamic systems involving non-linear and non-Gaussian properties. PF has been studied in various application areas including object tracking [Happe et al. 2011], robot localisation [Montemerlo et al. 2002], speech recognition [Vermaak et al. 2002] and air traffic management [Eele and Maciejowski 2011].

This work is supported in part by the European Union FP7 under grant agreement number 257906, 287804 and 318521, by UK EPSRC grant number EP/L00058X/1, EP/I012036/1 and EP/G066477/1, by Maxeler University Programme, by Xilinx, and by the Croucher Foundation.

Author's addresses: T.C.P. Chau, X. Niu and W. Luk, Department of Computing, Imperial College London; P.Y.K. Cheung, Department of Electrical and Electronic Engineering, Imperial College London; A. Eele and J. Maciejowski, Department of Engineering, University of Cambridge.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2014 ACM 1539-9087/2014/03-ART1 \$15.00

DOI: <http://dx.doi.org/10.1145/0000000.0000000>

PF keeps track of a large number of particles, each contains information about how a system would evolve. The underlying concept is to approximate a sequence of states by a collection of particles. Each particle is weighted to reflect the quality of an approximation. The more complex the problem, the larger the number of particles that are needed. One drawback of PF is its long execution times that limit its practical use.

This paper presents an efficient solution to PF. We derive an adaptive algorithm that adjusts its computation complexity at run time based on the quality of results. To map our algorithm to a heterogeneous reconfigurable system (HRS) consisting of multiple FPGAs and CPUs, we design a pipeline-friendly data structure to make effective use of the stream computing model. Moreover, we accelerate the algorithm with a data compression scheme and data control separation.

The key contributions of this paper include:

- (1) An adaptive PF algorithm which adapts the size of particle set at run-time. The algorithm is able to reduce computation workload while maintaining the quality of results.
- (2) Mapping the proposed algorithm to a scalable and reconfigurable system by following the stream computing model. A novel data structure is designed to take advantage of the architecture and to alleviate the data transfer bottleneck. The system uses the run-time reconfigurability of FPGA to switch between computation mode and low-power mode.
- (3) An implementation of a robot localisation application targeting the proposed system. Compared to a non-adaptive and non-reconfigurable implementation, the idle power of our proposed system is reduced by 25-34% and the overall energy consumption decreases by 17-33%. Our system with four FPGAs is up to 169 times faster than a single core CPU, 41 times faster than a 1U CPU server with 12 cores, and 3 times faster than a modelled four-GPU system.

2. BACKGROUND AND RELATED WORK

This section briefly outlines the PF algorithm. A more detailed description can be found in [Doucet et al. 2001]. PF estimates the state of a system by a sampling-based approximation of the state probability density function. The state of a system in time-step t is denoted by X_t . The control and observation are denoted by U_t and Y_t respectively. Three pieces of information about the system are known a-priori:

- $p(X_0)$ is the probability of the initial state of the system,
- $p(X_t|X_{t-1}, U_{t-1})$ is the state transition probability of the system's current state given a previous state and control information,
- $p(Y_t|X_t)$ is the observation model describing the likelihood of observing the measurement at the current state.

PF approximates the desired posterior probability $p(X_t|Y_{1:t})$ using a set of P particles $\{\chi_t^{(i)}\}_{i=1}^P$ with their associated weights $\{w^{(i)}\}_{i=1}^P$. X_0 and U_0 are initialised. This computation consists of three iterative steps.

- (1) **Sampling:** A new particle set $\{\tilde{\chi}_t^{(i)}\}_{i=1}^P$ is drawn from the distribution $p(X_t|X_{t-1}, U_{t-1})$, forming a prediction of the distribution of X_t .
- (2) **Importance weighting:** The likelihood $p(Y_t|\tilde{\chi}_t^{(i)})$ of each particle is calculated. The likelihood indicates whether the current measurement Y_t matches the predicted state $\{\tilde{\chi}_t^{(i)}\}_{i=1}^P$. Then each particle is assigned a weight $w^{(i)}$ with respect to the likelihood.
- (3) **Resampling:** Particles with higher weights are replicated and the number of particles with lower weights is reduced. With resampling, the particle set has a

smaller variance. The particle set is used in the next time-step to predict the posterior probability subsequently. The distribution of the resulting particles $\{\chi_t^{(i)}\}_{i=1}^P$ approximates $p(X_t|Y_{1:t})$.

The particles in PF are independent of each other. It means the algorithm can be accelerated using specialised hardware with massive parallelism and pipelining. In [Happe et al. 2011], an approach for PF on a hybrid CPU/FPGA platform is developed. Using a multi-threaded programming model, computation is switched between hardware and software during run-time to react to performance requirements. Resampling algorithms and architectures for distributed PFs are proposed in [Bolic et al. 2005].

Adaptive PFs have been proposed to improve performance or quality of state estimation by controlling the number of particles dynamically. Likelihood-based adaptation controls the number of particles such that the sum of weights exceeds a pre-specified threshold [Koller and Fratkinina 1998]. Kullback Leibler distance (KLD) sampling is proposed in [Fox 2003], which offers better quality results than likelihood-based approach. KLD sampling is improved in [Park et al. 2010] by adjusting the variance and gradient of data to generate particles near high likelihood regions. The above methods introduce data dependencies in the sampling and importance weighting steps, so they are difficult to be parallelised. An adaptive PF is proposed in [Bolic et al. 2002] that changes the number of particles dynamically based on estimation quality. In [Chau et al. 2012], adaptive PF is extended to a multi-processor system on FPGA. The number of particles and active processors change dynamically but the performance is limited by soft-core processors. In [Liu et al. 2007], a mechanism and a theoretical lower bound for adapting the sample size of particles are presented. Our previous work [Chau et al. 2013a] presents a hardware-friendly adaptive PF. The algorithm is mapped to an accelerator system which consists of an FPGA and a CPU. However, the system suffers from a large communication overhead when the particles are transferred between the FPGA and CPU. Moreover, the scalability of the adaptive PF algorithm to multiple FPGAs is not covered. In this paper, we extend our previous work to address the problems mentioned above.

3. ADAPTIVE PARTICLE FILTER

This section introduces an adaptive PF algorithm which changes the number of particles at each time-step. The algorithm is inspired by [Liu et al. 2007] and we transform it to a pipeline-friendly version for mapping to the stream computing architecture. This algorithm is shown in Algorithm 1 which consists of four stages.

3.1. Stage 1: Sampling and Importance Weighting (line 8 to 9)

At the initial time-step ($t = 0$), the maximum number of particles are used, i.e. $P_0 = P_{max}$. At the subsequent time-steps, the number of particles is denoted as P_t . Initially, the particle set $\{\chi_t^{(i)}\}_{i=1}^{P_t}$ is sampled to $\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}$. Then a weight from $\{w^i\}_{i=1}^{P_t}$ is assigned to each particle. As a result, $\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}$ and $\{w^{(i)}\}_{i=1}^{P_t}$ give an estimation of the next state.

During sampling and importance weighting, the computation of every particle is independent of each other. The mapping of computation to FPGAs will be described in Section 4.

3.2. Stage 2: Lower Bound Calculation (line 10)

This stage derives the smallest number of particles that are needed in the next time-step in order to bound the approximation error. The adaptive algorithm seeks a value

ALGORITHM 1: Adaptive PF algorithm

```

1:  $P_0 \leftarrow P_{max}$ 
2:  $\{\chi_0^{(i)}\}_{i=1}^{P_0} \leftarrow$  random set of particles
3:  $t = 1$ 
4: for each step  $t$  do
5:    $r = 0$ 
6:   while  $r \leq itl\_repeat$  do
7:     —On FPGAs—
8:     Sample a new state  $\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}$  from  $\{\chi_t^{(i)}\}_{i=1}^{P_t}$ 
9:     Calculate unnormalised importance weights  $\{\tilde{w}^{(i)}\}_{i=1}^{P_t}$  and accumulate the weights as  $w_{sum}$ 
10:    Calculate the lower bound of sample size  $\tilde{P}_{t+1}$  by Equation 1
11:    —On CPUs—
12:    Sort  $\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}$  in descending  $\{\tilde{w}^{(i)}\}_{i=1}^{P_t}$ 
13:    if  $\tilde{P}_{t+1} < P_t$  then
14:       $P_{t+1} = \max([\tilde{P}_{t+1}], P_t/2)$ 
15:      Set  $a = 2P_{t+1} - P_t$  and  $b = P_{t+1}$ 
16:      —Do the following loop in parallel—
17:      for  $i$  in  $P_t - P_{t+1}$  do
18:         $\tilde{\chi}_{t+1}^{(i)} = \frac{\chi_{t+1}^{(a)} \tilde{w}^{(a)} + \chi_{t+1}^{(b)} \tilde{w}^{(b)}}{\tilde{w}^{(a)} + \tilde{w}^{(b)}}$ 
19:         $\tilde{w}^{(i)} = \tilde{w}^{(a)} + \tilde{w}^{(b)}$ 
20:         $a = a + 1$  and  $b = b - 1$ 
21:      end for
22:      else if  $\tilde{P}_{t+1} \geq P_t$  then
23:         $a = 0$  and  $b = 0$ 
24:        for  $i$  in  $P_{t+1} - P_t$  do
25:          if  $\tilde{w}^{(a)} < \tilde{w}^{(a+1)}$  and  $a < P_{t+1}$  then
26:             $a = a + 1$ 
27:          end if
28:           $\tilde{\chi}_{t+1}^{(P_t+b)} = \tilde{\chi}_{t+1}^{(a)}/2$ 
29:           $\tilde{\chi}_{t+1}^{(a)} = \tilde{\chi}_{t+1}^{(a)}/2$ 
30:           $\tilde{w}^{(P_t+b)} = \tilde{w}^{(a)}/2$ 
31:           $\tilde{w}^{(a)} = \tilde{w}^{(a)}/2$ 
32:           $b = b + 1$ 
33:        end for
34:      end if
35:      Resample  $\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}$  to  $\{\chi_{t+1}^{(i)}\}_{i=1}^{P_{t+1}}$ 
36:       $r = r + 1$ 
37:    end while
38:  end for

```

which is less than or equal to P_{max} . This number, denoted as \tilde{P}_{t+1} , is referred to as the lower bound of sampling size. It is calculated by Equation 1 to 4.

$$\tilde{P}_{t+1} = \sigma^2 \cdot \frac{P_{max}}{Var(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t})} \quad (1)$$

$$\begin{aligned} \sigma^2 = & \sum_{i=1}^{P_t} \left(w^{(i)} \cdot \tilde{\chi}_{t+1}^{(i)} \right)^2 - 2 \cdot E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) \cdot \sum_{i=1}^{P_t} \left((w^{(i)})^2 \cdot \tilde{\chi}_{t+1}^{(i)} \right) \\ & + \left(E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) \right)^2 \cdot \sum_{i=1}^{P_t} (w^{(i)})^2 \end{aligned} \quad (2)$$

$$Var(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) = \sum_{i=1}^{P_t} \left(w^{(i)} \cdot (\tilde{\chi}_{t+1}^{(i)})^2 \right) - \left(E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) \right)^2 \quad (3)$$

$$E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) = \sum_{i=1}^{P_t} w^{(i)} \cdot \tilde{\chi}_{t+1}^{(i)} \quad (4)$$

As shown in Equation 2 to 4, $w^{(i)}$ is a normalised term. To calculate $w^{(i)}$, a traditional software-based approach is to iterate through the set of particles twice. The sum of weights w_{sum} and unnormalised weight $\tilde{w}^{(i)}$ are calculated in the first iteration. Then $w^{(i)}$ is obtained by dividing $\tilde{w}^{(i)}$ by w_{sum} in the second iteration. However, this method is inefficient for FPGA implementation. Since $2P_t$ cycles are needed to process P_t pieces of data, the throughput is reduced to 50%.

To fully utilise deep pipelines targeting an FPGA, we perform function transformation. Given $w^{(i)} = \frac{\tilde{w}^{(i)}}{w_{sum}}$, we extract w_{sum} out of Equation 2 to 4. By doing so, we obtain a transformed form as shown in Equations 5 to 7. w_{sum} and $\tilde{w}^{(i)}$ are computed simultaneously in two separate data paths. At the last clock cycle of the particle stream, σ^2 , $Var(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t})$ and $E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t})$ are obtained. The details of the FPGA kernel design will be explained in Section 4.

$$\begin{aligned} \sigma^2 = & \frac{1}{(w_{sum})^2} \cdot \left(\sum_{i=1}^{P_t} (\tilde{w}^{(i)} \cdot \tilde{\chi}_{t+1}^{(i)})^2 - 2 \cdot E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) \cdot \sum_{i=1}^{P_t} (\tilde{w}^{(i)})^2 \cdot \tilde{\chi}_{t+1}^{(i)} \right) \\ & + \left(E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) \right)^2 \cdot \sum_{i=1}^{P_t} (\tilde{w}^{(i)})^2 \end{aligned} \quad (5)$$

$$Var(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) = \frac{1}{w_{sum}} \cdot \sum_{i=1}^{P_t} (\tilde{w}^{(i)} \cdot (\tilde{\chi}_{t+1}^{(i)})^2) - \left(E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) \right)^2 \quad (6)$$

$$E(\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}) = \frac{1}{w_{sum}} \cdot \sum_{i=1}^{P_t} \tilde{w}^{(i)} \cdot \tilde{\chi}_{t+1}^{(i)} \quad (7)$$

3.3. Stage 3: Particle set size tuning (line 12 to 34)

The adaptive approach tunes the particle set size to fit the lower bound P_{t+1} . This stage is done on the CPUs because the operations involve non-sequential data access that cannot be mapped efficiently to FPGAs.

The particles are sorted in descending order according to their weights. As the new sample size can increase or decrease, there are two cases:

— Case I: Particle set reduction when $\tilde{P}_{t+1} < P_t$

The lower bound P_{t+1} is set to $\max(\lceil \tilde{P}_{t+1} \rceil, P_t/2)$. Since the new size is smaller than the old one, some particles are combined to form a smaller particle set. Figure 1 illustrates the idea of particle reduction. The first $2P_{t+1} - P_t$ particles with higher weights are kept and the remaining $2(P_t - P_{t+1})$ particles are combined in pairs. As a result, there are $P_t - P_{t+1}$ new particles injected to form the target particle set with P_{t+1} particles. We combine the particles deterministically to keep the statements in the loop independent of each other. As a result, loop unrolling is undertaken to execute the statements in parallel. The complexity of the loop is in $\mathcal{O}\left(\frac{P_t - P_{t+1}}{N_{parallel}}\right)$, where $N_{parallel}$ indicates the level of parallelism.

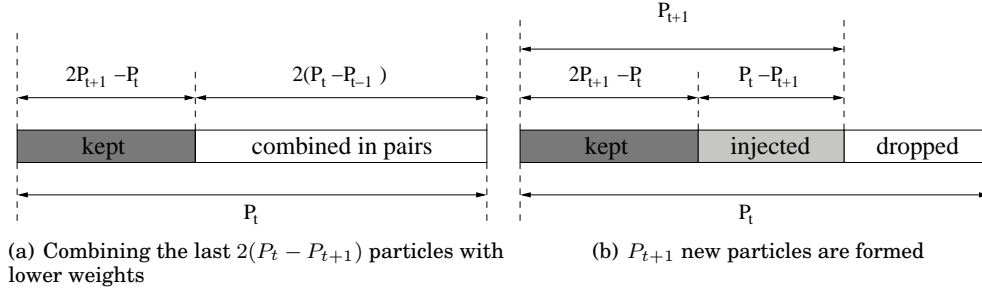


Fig. 1. Particle set reduction

— **Case II: Particle set expansion when $\tilde{P}_{t+1} \geq P_t$**

The lower bound P_{t+1} is set to \tilde{P}_{t+1} . Some particles are taken from the original set and are inserted to form a larger set. The particles with larger weight would have more descendants. As shown in line 22 to 34, the process requires picking the particle with the largest weight at each iteration of particle incision. Since the particle set is pre-sorted, the complexity of particle set expansion is $\mathcal{O}(P_{t+1} - P_t)$.

3.4. Stage 4: Resampling (line 35)

Resampling is performed to pick P_{t+1} particles from $\{\tilde{\chi}_{t+1}^{(i)}\}_{i=1}^{P_t}$ to form $\{\chi_{t+1}^{(i)}\}_{i=1}^{P_{t+1}}$. The process has a complexity of $\mathcal{O}(P_{t+1})$.

4. HETEROGENEOUS RECONFIGURABLE SYSTEM

This section describes the proposed heterogeneous reconfigurable system (HRS). It is scalable to cope with different FPGA devices and applications. HRS also takes advantage of the run-time reconfiguration feature for power and energy reduction.

4.1. Mapping adaptive PF to HRS

The system design of HRS is shown in Figure 2. A heterogeneous structure is employed to make use of multiple FPGAs and CPUs. FPGAs and CPUs communicate through high bandwidth buses. FPGAs are responsible for (1) sampling, (2) importance weighting, and (3) lower bound calculation. The data paths on the FPGAs are fully-pipelined. Each FPGA has its own on-board dynamic random-access memory (DRAM) to store the large amount of particle data. On the other hand, the CPUs gather all the particles from FPGAs to perform particle set size tuning and resampling.

4.2. FPGA Kernel Design

Sampling, importance weighting and lower bound calculation are the most computation intensive stages. In each time-step, these three stages are iterated for itl_repeat times. An FPGA kernel is designed to enable acceleration of them.

Figure 4 shows the components of the FPGA kernel. The kernel is fully pipelined to achieve one output per clock cycle. It can also be replicated as many times as FPGA resource allow and the replications can be split across multiple FPGA boards. The kernel takes three inputs from the CPUs or on-board DRAM: (1) states, (2) controls, and (3) seeds. Application specific parameters are stored in ROMs. Three building blocks correspond to the sampling, importance weighting and lower bound calculation stages as described in Section 3.

For sampling and importance weighting, the computation of each particle is independent of each other. Particles are fed to the FPGAs as a stream shown in Figure 3.

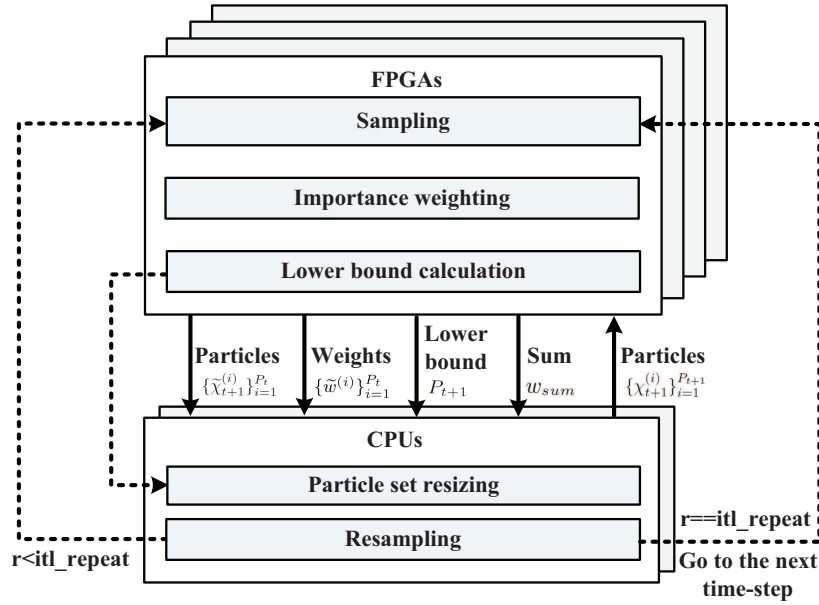


Fig. 2. Heterogeneous reconfigurable system (Solid lines: data paths; Dotted lines: control paths)

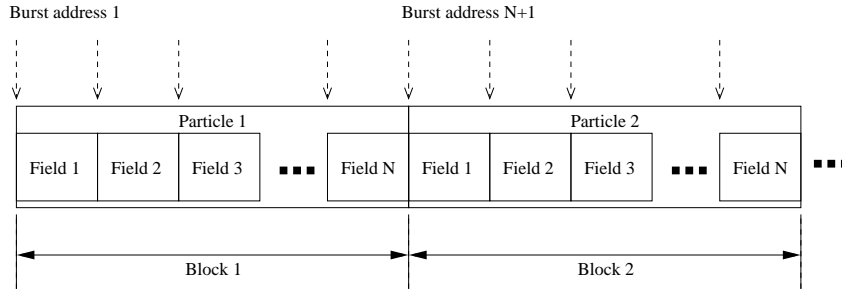


Fig. 3. A particle stream

Each block of the particle stream consists of a number of data fields which store information of a particle. The number of data fields is application dependent. In every clock cycle, one piece of data is transferred from the onboard memory to an FPGA data path. Each FPGA data path has a long pipeline where each stage is filled with a piece of data, and therefore many particles are processed simultaneously. Fixed-point data representation is customised at each pipeline stage to reduce the resource usage.

Meanwhile, the accumulation of w_{sum} introduces a feedback loop. A new weight comes along every cycle which is more quickly than the floating-point unit to perform addition of the previous weight. In order to achieve one result per clock cycle, fixed-point data-path is implemented while ensuring no overflow or underflow occurs.

4.3. Timing model for run time reconfiguration

We derive a model to analyse the computation time of HRS. The model helps us to design a configuration schedule that satisfies the real-time requirement and, if neces-

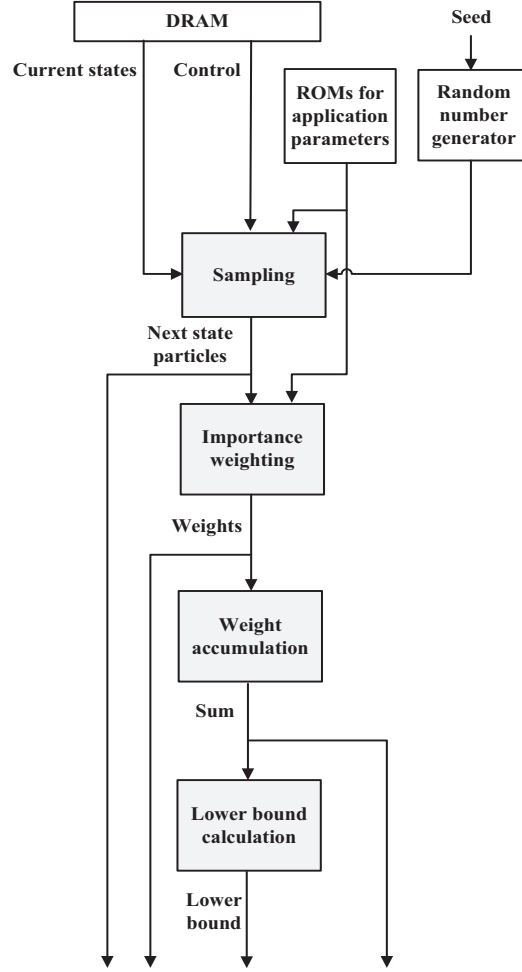


Fig. 4. FPGA kernel design

sary, amend the application's specification. The model will be validated by experiments in Section 6.

The computation time (T_{comp}) of HRS consists of three components: (1) Data path time $T_{datapath}$, (2) CPU time T_{CPU} , and (3) Data transfer time T_{tran} . The sampling, importance weighting and resampling processes are repeated for itl_repeat times in every time-step.

$$T_{comp} = itl_repeat \cdot (T_{datapath} + T_{CPU} + T_{tran}) \quad (8)$$

Data path time, $T_{datapath}$, denotes the time spent on the FPGAs. P_t denotes the number of particles at the current time-step and f_{FPGA} denotes the clock frequency of the FPGAs. L is the length of the pipeline. $N_{datapath}$ denotes the number of data paths on one FPGA board. N_{FPGA} is the number of FPGA boards in the system.

$$T_{datapath} = \left(\frac{P_t}{f_{FPGA} \cdot N_{datapath}} + L - 1 \right) \frac{1}{N_{FPGA}} \quad (9)$$

CPU time, T_{CPU} , denotes the time spent on the CPUs. The clock frequency and number of threads of the CPUs are represented by f_{CPU} and N_{thread} respectively. par is an application-specific parameter in the range of $[0, 1]$ which represents the ratio of CPU instructions that are parallelisable, and α is a scaling constant derived empirically.

$$T_{CPU} = \alpha \cdot \frac{P_t}{f_{CPU}} \cdot \left(1 - par + \frac{par}{N_{thread}} \right) \quad (10)$$

Data transfer time, T_{tran} , denotes the time of moving a particle stream between the FPGAs and the CPUs. df is the number of data fields of a particle. For example, if a particle contains the information of coordinates (x, y) and heading h , $df = 3$. Given that the constant 1 represents the weight and the constant 2 accounts for the movement of data in and out of the FPGAs, and bw_{data} is the bit-width of one data field, the expression $(2 \cdot df + 1) \cdot bw_{data}$ is regarded as the size of a particle.

f_{bus} is the clock frequency of the bus connecting the CPUs to FPGAs and $lane$ is the number of bus lanes connected to one FPGA. Since many buses, such as the PCI Express Bus, encode data during transfer, the effective data are denoted by eff (in PCI Express Gen2 the value is 8/10). In our previous work [Chau et al. 2013a], the data transfer time has a significant performance impact on HRS. To reduced the data transfer overhead, we introduce a data compression technique that will be described in Section 5.

$$T_{tran} = \frac{(2 \cdot df + 1) \cdot bw_{data} \cdot P_t}{f_{bus} \cdot lane \cdot eff \cdot N_{FPGA}} \quad (11)$$

In real-time applications, each time-step is fixed and is known as the real-time bound T_{rt} . The derived model helps system designers to ensure that the computation time T_{comp} is shorter than T_{rt} . An idle time T_{idle} is introduced to represent the time gap between the computation time and real-time bound.

$$T_{idle} = T_{rt} - T_{comp} \quad (12)$$

Figure 5(a) illustrates the power consumption of an HRS without run-time reconfiguration. It shows that the FPGAs are still drawing power after the computation finishes. By exploiting run-time reconfiguration as shown in Figure 5(b), the FPGAs are loaded with a low-power configuration during the idle period. Such configuration minimises the amount of active resources and clock frequency. Equation 13 describes the sleep time when the FPGAs are idle and being loaded with the low-power configuration. If the sleep time is positive, reconfiguration would be helpful in these situations.

$$T_{sleep} = T_{idle} - T_{config} \quad (13)$$

Configuration time, T_{config} , denotes the time needed to download a configuration bit-stream to the FPGAs. $size_{bs}$ represents the size of bitstream in bits. f_{config} is the configuration clock frequency in Hz and bw_{config} is the width of the configuration port.

$$T_{config} = \frac{size_{bs}}{f_{config} \cdot bw_{config}} \quad (14)$$

5. OPTIMISING TRANSFER OF PARTICLE STREAM

In Section 4, the data transfer time depends on the number of particles and the bus bandwidth between the CPUs and FPGAs. It can be a major performance bottleneck

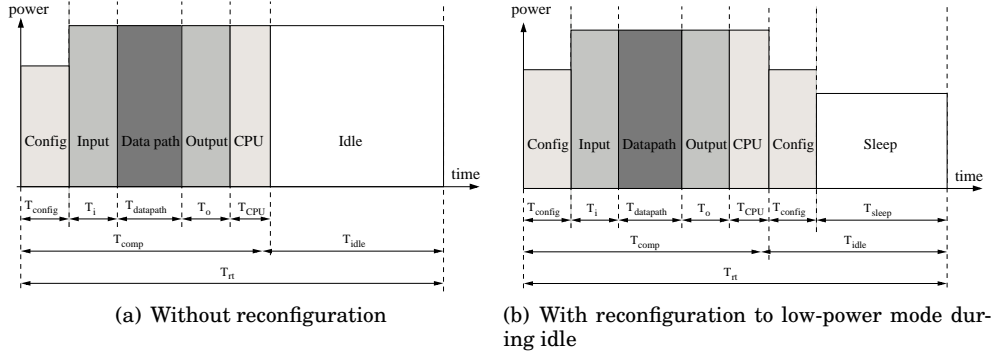


Fig. 5. Power consumption of the HRS over time

as depicted in [Chau et al. 2013a]. Refer to Figure 6(a), each block stores the data of a particle. When the CPUs finish processing, all data are transferred from the CPUs to the FPGAs. The data transfer time cannot be reduced by implementing more FPGA data paths or increasing the FPGAs' clock frequency because the bottleneck is at the bus connecting the CPUs and FPGAs.

To improve the data transfer performance, we design a data structure which facilitates compression of particles. The idea comes from an observation of the resampling process - some particles are eliminated and the vacancies are filled by replicating non-eliminated particles. Replication means data redundancy exists. For example, in the original data structure shown in Figure 6(a), particle 1 has three replicates and particle 2 is eliminated, therefore, particle 1 is stored and transferred for three times.

By using the data structure in Figure 6(b), data redundancy is eliminated by storing every particle once. Each particle is also transferred once. As a result, the data transfer time and memory space are reduced.

An HRS often contains DRAM which transfers data in burst in order to maximise the memory bandwidth. This works fine with the original data structure where the data are organised as a sequence from the lower address space to the upper. However, using the new data structure, the data access pattern is not sequential anymore, the address can go back and forth. The DRAM controller needs to be modified so that the transfer throughput would not be affected by the change of data access pattern. As illustrated in Figure 6(b), a tag sequence is used to indicate the address of the next block. For example, after reading the data of particle 1, the burst address is at N . If the tag is one, the next burst address will point to the address of the next block at $N + 1$. Otherwise, the burst address will point to the start address of the current block (which is 1). The data are still addressed in burst so the performance is not degraded.

The data transfer time with compression is shown below. Rep is the average number of replication of the particles, and therefore the size of the resampled particle stream is reduced by a ratio of Rep . The range of Rep is from 1 to P_t , depending on the distribution of particles after the resampling process. The effect of Rep on data transfer time will be evaluated in the next section.

$$T_{tran} = \frac{(\frac{df}{Rep} + df + 1) \cdot bw_{data} \cdot P_t}{f_{bus} \cdot lane \cdot eff \cdot N_{FPGA}} \quad (15)$$

6. EXPERIMENTAL RESULTS

To evaluate the performance of the HRS and make comparison with the other systems, we implement an application which uses PF for localisation and tracking of mo-

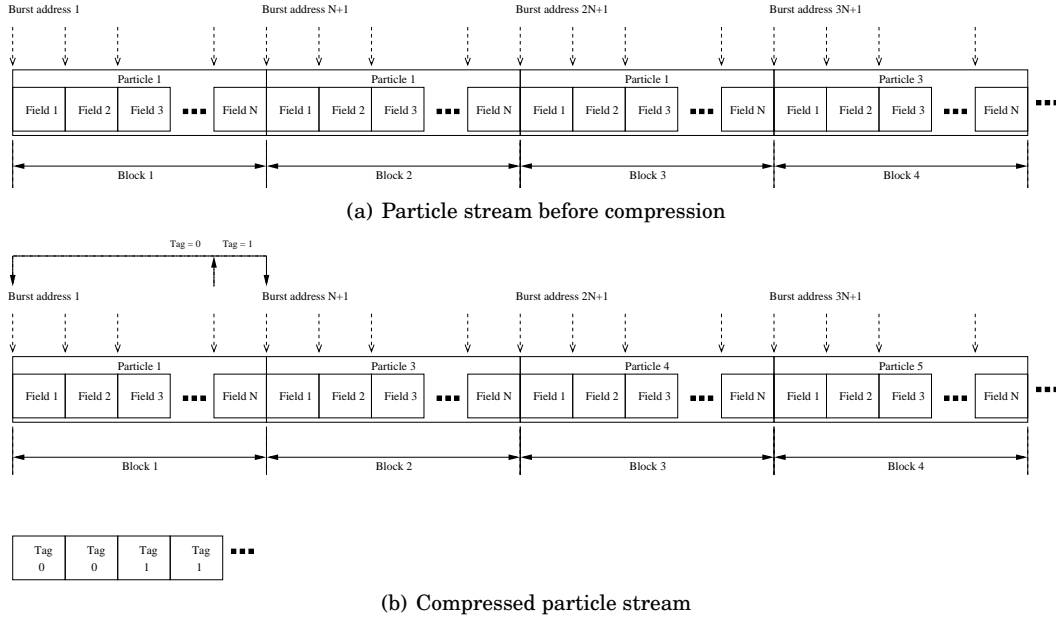


Fig. 6. After the resampling process, some particles are eliminated and the remaining particles are replicated. Data compression is applied so that every particle is stored and transferred once only.

bile robot. The application is proposed in [Montemerlo et al. 2002] to track location of moving objects conditioned upon robot poses over time. Given a priori learned map, a robot receives sensor values and moves at regular time intervals. Meanwhile, M moving objects are tracked by the robot. The states of the robot and objects at time t are represented by a state vector X_t :

$$X_t = \{R_t, H_{t,1}, H_{t,2}, \dots, H_{t,M}\} \quad (16)$$

R_t denotes the robot's pose at time t , and $H_{t,1}, H_{t,2}, \dots, H_{t,M}$ denote the locations of the M objects at the same time.

The following equation is used to represent the posterior of the robot's location:

$$p(X_t|Y_t, U_t) = p(R_t|Y_t, U_t) \prod_{m=1}^M p(H_{t,m}|R_t, Y_t, U_t) \quad (17)$$

Y_t is the sensor measurement and U_t is the control of the robot at time t . The robot path posterior $p(R_t|Y_t, U_t)$ is represented by a set of robot-particles. The distribution of an object's location $p(H_{t,m}|R_t, Y_t, U_t)$ is represented by a set of object-particles, where each object-particle set is attached to one particular robot-particle. In other words, if there are P_r robot-particles representing the posterior over robot path, there are P_r object-particle sets, each has P_h particles.

In the application, the area of the map is 12m by 18m. The robot makes a movement of 0.5m every five seconds, i.e. $T_{rt} = 5$. The robot can track eight moving objects at the same time. A maximum of 8192 particles are used for robot-tracking and each robot-particle is associated with 1024 object-particles. Therefore, the maximum number of data path cycles is $8 \cdot 8192 \cdot 1024 = 67,108,864$. Each particle being streamed into the FPGAs contains coordinates (x, y) and heading h which are represented by three single precision floating-point numbers. For the particle being streamed out of the FPGAs, it

also contains a weight in addition to the coordinates. From Equation 11, the size of a particle is $(2 \cdot 3 + 1) \cdot 32 \text{ bits} = 224 \text{ bits}$.

6.1. System Settings

HRS: Two reconfigurable accelerator systems from Maxeler Technologies are used. The system is developed using MaxCompiler, which adopts a stream computing model.

- *MaxWorkstation* is a microATX form factor system which is equipped with one Xilinx Virtex-6 XC6VVSX475T FPGA. The FPGA has 297,600 lookup tables (LUTs), 595,200 flip-flops (FFs), 2,016 digital signal processors (DSPs) and 1,064 block RAMs. The FPGA board is connected to an Intel i7-870 CPU (4 physical cores, 8 threads in total, clocked at 2.93 GHz) via a PCI Express Gen2 x8 bus. The maximum bandwidth of the PCI Express bus is 2 GB/s according to the specification provided by Maxeler Technologies.
- *MPC-C500* is a 1U server accommodating four FPGA boards, each of which has a Xilinx Virtex-6 XC6VVSX475T FPGA. Each FPGA board is connected to two Intel Xeon X5650 CPUs (12 physical cores, 24 threads in total, clocked at 2.66 GHz) via a PCI Express Gen2 x8 bus.

To support run-time reconfigurability, there are two FPGA configurations:

- *Sampling and importance weighting configuration* is clocked at 100 MHz. Two data paths are implemented on one FPGA to process particles in parallel. The total resource usage is 231,922 LUTs (78%), 338,376 FFs (56%), 1,934 DSPs (96%) and 514 block RAMs (48%).
- *Low-power configuration* is clocked at 10 MHz, with 5,962 LUTs (2%), 6,943 FFs (1%) and 12 block RAMs (1%). It uses minimal resources just to maintain communication between the FPGAs and CPUs.

CPU: The CPU performance results are obtained from a 1U server that hosts two Intel Xeon X5650 CPUs. Each CPU is clocked at 2.66 GHz. The program is written in C language and optimised by Intel Compiler with SSE4.2 and flag *-fast* enabled. OpenMP is used to utilise all the processor cores.

GPU: An NVIDIA Tesla C2070 GPU is hosted inside a 4U server. It has 448 cores running at 1.15 GHz and has a peak performance by 1288 GFlops. The program is written in C for CUDA and optimised to use all the cores available. To get more comprehensive results for comparison, we also estimate the performance of multiple GPUs. The estimation is based on the fact that the first three stages (sampling, importance weighting, lower bound calculation) can be evenly distributed to every GPU and be computed independently, so the data path and data transfer speedup scales linearly with the number of GPUs. On the other hand, the last two stages (particle set resizing, resampling) are computed on the CPU no matter how many GPUs are used, therefore, the CPU time does not scale with the number of GPUs.

6.2. Adaptive PF versus Non-adaptive PF

The comparison of adaptive and non-adaptive PF is shown in Table I. Both model estimation and experimental results are listed. Initially, the maximum number of particles are instantiated for global localisation. For the non-adaptive scheme, the particle set size does not change. The total computation time estimated and measured are 1.328 seconds and 1.885 seconds, respectively. The difference is due to the difference between the effective and maximum bandwidth of the PCI Express bus.

For the adaptive scheme, the number of particles varies from 573k to 67M, and the computation time scales linearly with the number of particles. From Table I, both the model and experiment show 99% reduction in computation time.

Table I. Comparison of adaptive and non-adaptive PF on HRS (MaxWorkstation with one FPGA, no data compression is applied)

	Non-adaptive PF		Adaptive PF	
	Model	Exp.	Model	Exp.
No. of particles	67M		573k	
Data path time $T_{datapath}$ (s)	0.336	0.336	0.003	0.003
CPU time T_{CPU} (s)	0.117	0.117	0.001	0.001
Data time T_{tran} (s)	0.875	1.432	0.007	0.012
Total comp. time T_{comp} (s)	1.328	1.885	0.011	0.016
Comp. speedup (higher is better)	1x	1x	120.7x	117.8x

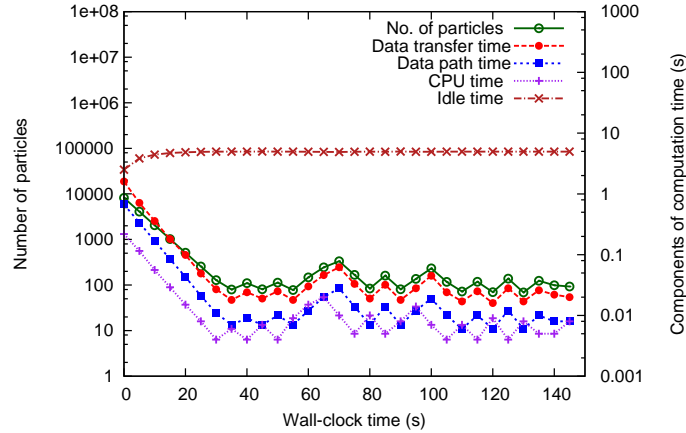


Fig. 7. Number of particles and components of total computation time versus wall-clock time

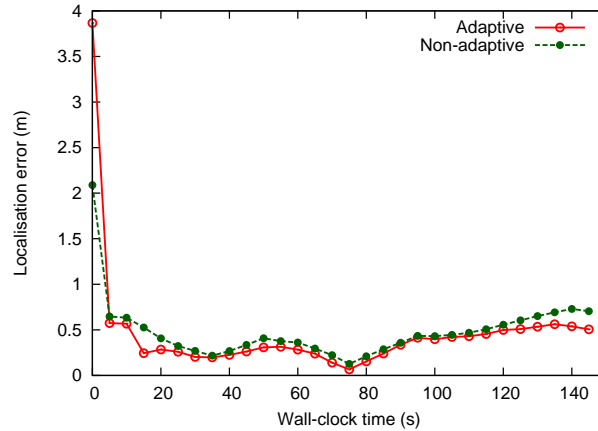


Fig. 8. Localisation error versus wall-clock time

Figure 7 shows how the number of particles and the components of total computation time vary over the wall-clock time (passage of time from the start to the completion of the application). Although the number of particles is reduced in the proposed design, the results in Figure 8 show that the localisation error is not adversely affected. The error is the highest during initial global localisation and it is reduced when the robot moves.

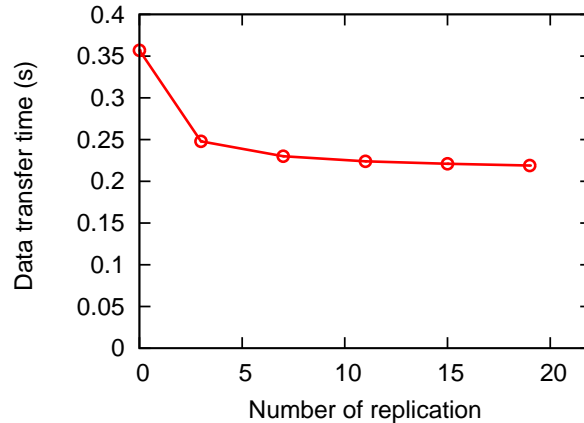


Fig. 9. Effect on the data transfer time by particle stream compression

6.3. Data Compression

Figure 9 shows the reduction in data transfer time after applying data compression. A higher number of replications means a lower data transfer time. The data transfer time has a lower bound of 0.212 seconds because the data from the FPGAs to the CPUs are not compressible. Only the particle stream after the resampling process is compressed when it is transferred from the CPUs to the FPGAs.

6.4. Performance comparison of HRS, CPUs and GPUs

Table II shows the performance comparison of the CPUs, GPUs and HRS.

Data path time: Considering the time spent on the data paths only, HRS is up to 328 times faster than a single-core CPU and 76 times faster than a 12-core CPU system with 24 threads. In addition, it is 12 times and 3 times faster than one GPU and four GPUs, respectively.

Data transfer time: The data transfer time of HRS is shown in three rows. The first row shows the situation when the PCI Express bandwidth is 2 GB/s. The second row shows the performance when PCI Express gen3 x8 (7.88 GB/s) is used such that the bandwidth is comparable with that of the GPU system. When multiple FPGA boards are used, the data transfer time decreases because multiple PCI Express buses are utilised simultaneously. The third row shows the performance when data compression is applied and it is assumed that each particle is replicated for 20 times in average.

CPU time: The CPU time of HRS is shorter than that of the CPU and GPU systems because part of the resampling process of object-particles is performed on the FPGA using Independent Metropolis-Hastings (IMH) resampling algorithm [Miao et al. 2011]. IMH resampling algorithm is optimised for the deep pipeline architecture where each particle occupies a single stage of the pipeline. On the CPUs and GPU, the computation of the particles are shared by threads and therefore IMH resampling algorithm is not applicable.

Total computation time: Considering the overall system performance, HRS is up to 169 times faster than a single-core CPU, 41 times faster than a 12-core CPU system. In addition, it is 9 times faster than one GPU, and 3 times faster than four GPUs. Notice that the CPUs violate the real-time constraint of 5 seconds.

Power and energy consumption: In real-time applications, we are interested in the energy consumption per time-step. Figure 10 shows the power consumption of HRS, CPUs and GPU over a period of 10 seconds (2 time-steps). The system power is

Table II. Performance comparison of HRS, CPUs and GPU

	CPU(1) ^a	CPU(2) ^a	GPU(1) ^b	GPU(2) ^b	GPU(3) ^b	HRS(1) ^c	HRS(2) ^d	HRS(3) ^d
Clock freq. (MHz)	2660	2660	1150	1150	1150	100	100	100
Precision	single	single	single	single	single	single + custom	single + custom	single + custom
Level of parallelism	1	24	448	896	1792	2+8 ^e	4+24 ^e	8+24 ^e
Data path time (s)	27.530	6.363	1.000	0.500	0.250	0.336	0.168	0.084
Data path speedup	1x	4.3x	27.5x	55.1x	110.1x	81.9x	163.9x	327.7x
Data tran. time (s)	0	0	0.360	0.180	0.090	1.432 ^f	0.716 ^f	0.358 ^f
CPU time (s)	0.420	0.334	0.117	0.117	0.117	0.363 ^g	0.182 ^g	0.091 ^g
Total comp. time (s)	27.95	6.697	1.477	0.797	0.457	0.223 ^h	0.111 ^h	0.056 ^h
Overall speedup	1x	4.2x	18.9x	35.1x	61.2x	47.5x	91.9x	169.4x
Comp. power (W)	183	279	287	424	698	145	420	480
Comp. power eff.	1x	0.7x	0.6x	0.4x	0.3x	1.3x	0.4x	0.4x
Idle power (W)	133	133	208	266	382	95	360	360
Idle power eff.	1x	1x	0.6x	0.5x	0.4x	1.4x	0.4x	0.4x
Energy. (J) ⁱ	677/5115	673/1868	1041/1157	1331/1456	1911/2054	489/595	1896/1914	1994/2012
Energy eff.	1x	1x/2.7x	0.7x/4.4x	0.5x/3.5x	0.4x/2.5x	1.4x/8.6x	0.4x/2.7x	0.3x/2.5x

a 2 Intel Xeon X5650 CPUs @2.66 GHz (12 cores supporting 24 threads).

b 1/2/4 NVIDIA Tesla C2070 GPUs and 1 Intel Core i7-950 CPU @3.07 GHz (4 cores supporting 8 threads).

c 1 Xilinx XC6VVSX475T FPGA and 1 Intel Core i7-870 CPU @2.93 GHz (4 cores supporting 8 threads).

d 4 Xilinx XC6VVSX475T FPGAs and 2 Intel Xeon X5650 CPUs @2.66 GHz (12 cores supporting 24 threads).

e Number of FPGA data paths and number of CPU threads.

f Each FPGA communicates with CPUs via a PCI Express bus with 2 GB/s bandwidth.

g Each FPGA communicates with CPUs via a PCI Express Gen3 x8 bus with 7.88 GB/s bandwidth.

h Each FPGA communicates with CPUs via a PCI Express Gen3 x8 bus with data compression.

i Cases for 573k and 67M particles in a 5-second interval.

measured using a power meter which is connected directly between the power source and the system. All the curves of HRS show peaks when HRS is at the computation mode and troughs when it is at the low power mode. The power during the configuration period lies between the two modes. On the HRS with one FPGA, run-time reconfiguration reduces the idle power consumption by 34% from 145W to 95W. In other words, over a 5-second time-step, the energy consumption is reduced by up to 33%. On the HRS with four FPGAs, the idle power consumption is reduced by 25% from 480W to 360W, and hence the energy consumption decreased by up to 17%.

The run-time reconfiguration methodology is not limited to the Maxeler systems, it can be applied to other FPGA platforms. The resource management software of our system (MaxelerOS) simplifies the effort of performing run-time reconfiguration, and hence we can focus on studying the impact of run-time reconfiguration on energy saving.

To identify the speed and energy trade-off, we produce a graph as shown in Figure 11. Each data point represents the computation time versus energy consumption of a system setting. Among all the systems, the HRS with one FPGA has the computation speed that satisfies the real-time requirement, while at the same time consumes the smallest amount of energy. All the configurations of CPU system cannot meet the real-time requirement. HRS(3), the HRS with four FPGAs, is the fastest among all the systems in comparison, therefore it is able to handle larger problems and more complex applications.

7. CONCLUSION

This paper presents an approach for accelerating adaptive particle filter for real-time applications. The proposed heterogeneous reconfigurable system demonstrates a significant reduction in power and energy consumption compared with CPU and GPU.

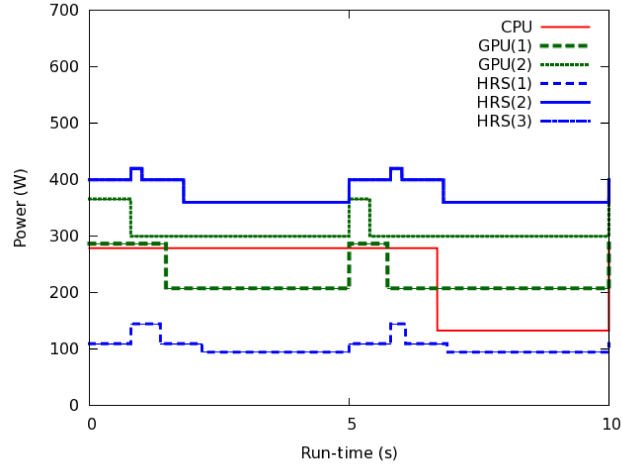


Fig. 10. Power consumption of HRS, CPU and GPU in one time-step, notice that the computation time of the CPU system exceeds the 5-second real-time requirement (The lines of HRS(2) and HRS(3) overlap)

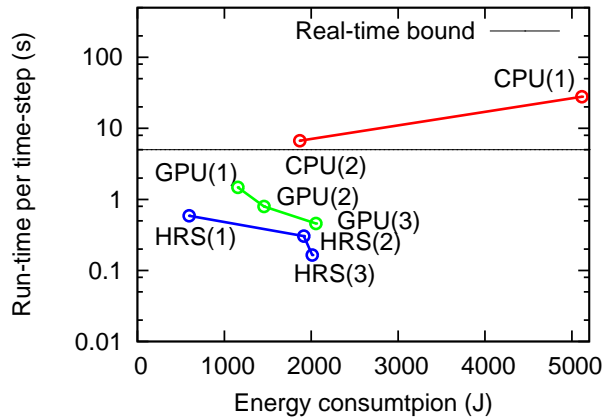


Fig. 11. Run-time vs. energy consumption of HRS, CPUs and GPUs (5-second time-step, 67M particles; Refer to Table II for system settings)

The adaptive algorithm reduces computation time while maintaining the quality of results. The approach is scalable to systems with multiple FPGAs. A data compression technique is used to mitigate the data transfer overhead between the FPGAs and CPUs.

In the future, heterogeneous reconfigurable systems will be developed for various particle filters that are more compute-intensive and have more stringent real-time requirements than the ones described above. Air traffic management [Chau et al. 2013b] and traffic estimation [Mihaylova et al. 2007] are example applications that can substantially benefit from the proposed approach in meeting current and future requirements. Further work will also be required to automate the optimisation of designs targeting heterogeneous reconfigurable systems.

ACKNOWLEDGMENTS

This work is supported in part by the European Union Seventh Framework Programme under grant agreement number 257906, 287804 and 318521, by UK EPSRC grant number EP/L00058X/1, EP/I012036/1 and EP/G066477/1, by Maxeler University Programme, by Xilinx, and by the Croucher Foundation. The authors thank Oliver Pell at Maxeler Technologies for comments on the paper.

REFERENCES

- Miodrag Bolic, Petar M. Djuric, and Sangjin Hong. 2005. Resampling algorithms and architectures for distributed particle filters. *IEEE Trans. Signal Processing* 53, 7 (2005), 2442–2450.
- Miodrag Bolic, Sangjin Hong, and Petar M. Djuric. 2002. Performance and complexity analysis of adaptive particle filtering for tracking applications. In *Proc. Asilomar Conf. Signals, Systems, and Computers*, Vol. 1. 853–857.
- Thomas C.P. Chau, Wayne Luk, Peter Y.K. Cheung, Alison Eele, and Jan Maciejowski. 2012. Adaptive Sequential Monte Carlo Approach for Real-Time Applications. In *Proc. Int. Conf. Field Programmable Logic and Applications*. 527–530.
- Thomas C.P. Chau, Xinyu Niu, Alison Eele, Wayne Luk, Peter Y.K. Cheung, and Jan Maciejowski. 2013a. Heterogeneous Reconfigurable System for Adaptive particle Filters in Real-Time Applications. In *Proc. Int. Symp. Applied Reconfigurable Computing*. 1–12.
- Thomas C.P. Chau, James S. Targett, Marlon Wijeyasinghe, Wayne Luk, Peter Y.K. Cheung, Benjamin Cope, Alison Eele, and J.M. Maciejowski. 2013b. Accelerating Sequential Monte Carlo Method for Real-time Air Traffic Management. In *Proc. Int. Symp. Highly Efficient Accelerators and Reconfigurable Technologies*.
- Arnaud Doucet, Nando de Freitas, and Neil Gordon. 2001. *Sequential Monte Carlo methods in practice*. Springer.
- A. Eele and J.M. Maciejowski. 2011. Comparison of Stochastic Optimisation Methods for Control in Air Traffic Management. In *Proc. IFAC World Congress*.
- Dieter Fox. 2003. Adapting the sample size in particle filters through KLD-sampling. *Int. Trans. Robotics* 22, 12 (2003), 985–1003.
- Markus Happe, Enno Lübbers, and Marco Platzner. 2011. A self-adaptive heterogeneous multi-core architecture for embedded real-time video object tracking. *Journal Real-Time Image Processing* (2011), 1–16.
- Daphne Koller and Raya Fratkina. 1998. Using learning for approximation in stochastic processes. In *Proc. Int. Conf. Machine Learning*. 287–295.
- Zhibin Liu, Zongying Shi, Mingguo Zhao, and Wenli Xu. 2007. Mobile robots global localization using adaptive dynamic clustered particle filters. In *Proc. Int. Conf. Intelligent Robots and Systems*. 1059–1064.
- Lifeng Miao, Jun Jason Zhang, Chaitali Chakrabarti, and Antonia Papandreou-Suppappola. 2011. Algorithm and Parallel Implementation of Particle Filtering and its Use in Waveform-Agile Sensing. *J. Signal Process. Syst.* 65, 2 (2011), 211–227.
- Lyudmila Mihaylova, Ren. Boel, and Andreas Hegyi. 2007. Freeway traffic estimation within particle filtering framework. *Automatica* 43, 2 (2007), 290–300.
- Michael Montemerlo, Sebastian Thrun, and William Whittaker. 2002. Conditional particle filters for simultaneous mobile robot localization and people-tracking. In *Proc. Int. Conf. Robotics and Automation*. 695–701.
- Sang-Hyuk Park, Young-Joong Kim, and Myo-Taeg Lim. 2010. Novel adaptive particle filter using adjusted variance and its application. *Int. Journal Control, Automation and Systems* 8, 4 (2010), 801–807.
- Jaco Vermaak, Christophe Andrieu, Arnaud Doucet, and Simon John Godsill. 2002. Particle methods for Bayesian modeling and enhancement of speech signals. *IEEE Trans. Speech and Audio Processing* 10, 3 (2002), 173–185.

Received June 2013; revised February 2014; accepted March 2014